

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-344282

(43)Date of publication of application : 14.12.2001

(51)Int.Cl.

G06F 17/30

(21)Application number : 2001-064336

(71)Applicant : HITACHI LTD

(22)Date of filing : 28.02.1991

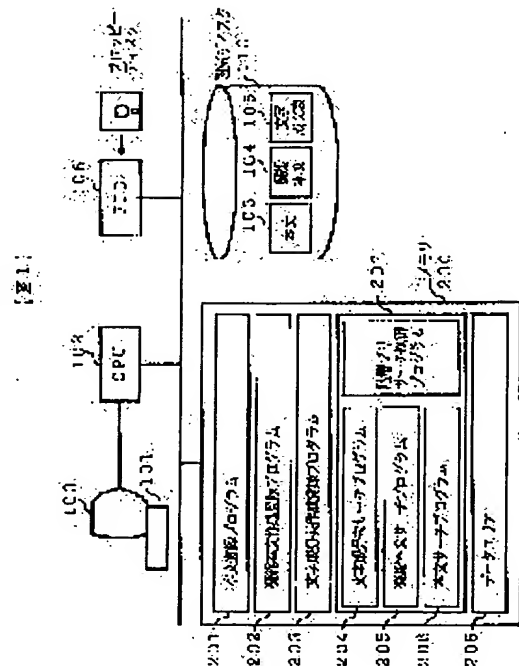
(72)Inventor : HATAKEYAMA ATSUSHI  
FUJISAWA HIROMICHI  
KATO KANJI  
KAWAGUCHI HISAMITSU  
MINEGISHI NAOKI

## (54) METHOD AND DEVICE FOR DOCUMENT RETRIEVAL

## (57)Abstract:

**PROBLEM TO BE SOLVED:** To solve such problem that the conventional full text search system requires much processing time because an entire text file is scanned from the beginning character by character and this system requires much time when being applied to a large-scale database.

**SOLUTION:** When a document is registered in a document database, a text character string of the registered document is divided for every character type such as Hiragana (cursive form of Japanese syllabary) or KANJI (Chinese character), and the inclusion relations between divided partial character strings are examined to generate a condensed text consisting of a set of partial character strings where character strings included in other character strings are excluded, and a character component table is generated where characters appearing in the condensed text are registered without duplication, and the condensed text and the character component table are registered in the document database besides the text of the registration object document. At the time of retrieval, the character component table is referred to extract a document including characters of a designated keyword, and next, the condensed text of the extracted document is referred to extract only documents corresponding to the condensed text including the partial character string of the designated keywords, and texts of extracted documents are referred to extract only a document meeting a retrieval condition given between keywords.



## LEGAL STATUS

[Date of request for examination] 08.03.2001

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3303881

[Date of registration] 10.05.2002

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開2001-344282

(P2001-344282A)

(43)公開日 平成13年12月14日(2001.12.14)

(51) Int.Cl.<sup>7</sup>

識別記号

FI

テーマート\* (参考)

G O 6 F 17/30

4 1 4

G O 6 F 17/30

4 1 4 B      5 B 0 7 5

170

170A

220

220A

審査請求 有 請求項の数 4 O.L (全 27 頁)

(21)出願番号 特願2001-64336(P2001-64336)

(62) 分割の表示 特願平3-58311の分割

(22)出願日 平成3年2月28日(1991.2.28)

(71)出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72)発明者 畠山 敦

東京都国分寺市東恋ヶ窪1丁目280番地

株式会社日立製作所中央研究所内

(72)発明者 藤澤 浩道

東京都国分寺市東恋ヶ窪1丁目280番地

株式会社日立製作所中央研究所内

(74)代理人 100075096

弁理士 作田 康夫

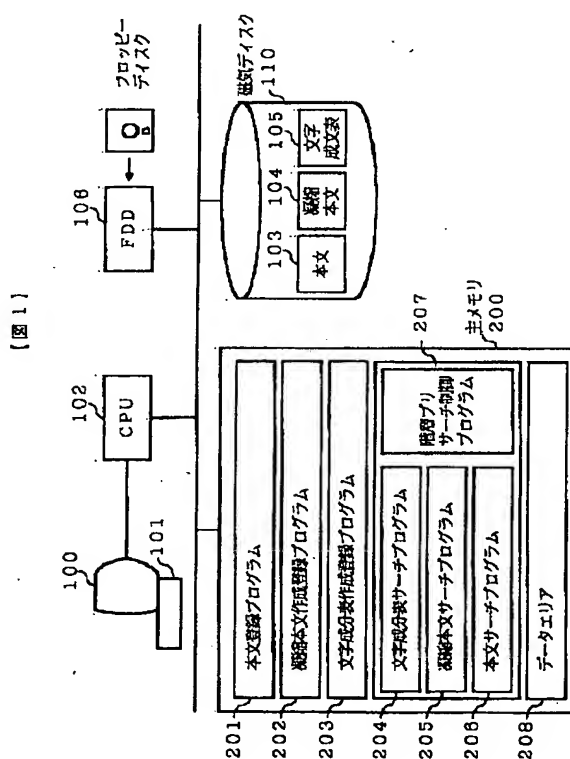
[最終頁に続く](#)

(54) 【発明の名称】 文書検索方法および装置

(57) 【要約】

【課題】 従来のフルテキストサーチ方式には、テキストファイル全体を先頭から一文字ずつ走査するために処理時間が掛かり、大規模なデータベースに適用すると時間が掛かるという課題があった。

【解決手段】 文書データベースに文書を登録する際、登録文書の本文文字列をひらがな、漢字等の文字種ごとに分割し、分割した各部分文字列の間で相互に文字列の包含関係を調べ、他の文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成し、凝縮本文中に現れる文字を重複なく登録した文字成分表を作成し、登録対象文書の本文のほかに凝縮本文と文字成分表を合わせて文書データベースへ登録する。検索時には指定キーワードの文字を含む文書を、文字成分表を参照して抽出し、次に抽出された文書の凝縮本文を参照して指定キーワードの部分文字列を含む凝縮本文に対応する文書のみを抽出し、該抽出文書の本文を参照してキーワード間に付与された検索条件を満たすもののみを抽出する。



**【特許請求の範囲】**

【請求項1】文書情報を文字コードデータとして蓄積した文書データベースを対象として、検索者が指定したキーワードを含む文書を検索する文書検索装置において、該文書データベースに文書を登録する際、該登録文書の本文文字列をひらがな、漢字、及び英数字等の文字種ごとに分割し、分割した各部分文字列の間で相互に文字列の包含関係を調べ、他の文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成する手段と、該凝縮本文中に現れる文字を重複なく登録した文字成分表を作成する手段と、登録対象文書の本文のほかに凝縮本文と文字成分表を合わせて文書データベースへ格納する手段と、検索時に検索者が指定したキーワードを構成する全ての種類の文字を含む文書を、文字成分表を参照して抽出する文字成分表サーチ手段と、該文字成分表サーチで抽出された文書の件数を検査する手段と、該検査の結果件数が所定数以上の件数に達した場合は、凝縮本文を全件読み出して、検索者が指定したキーワードを構成する部分文字列を含む凝縮本文に対応する文書を抽出し、前記検査の結果件数が所定数以下の場合は、前記文字成分表サーチで抽出された文書の凝縮本文を参照して、検索者が指定したキーワードを構成する部分文字列を含む凝縮本文に対応する文書を抽出する凝縮本文サーチ手段と、抽出された文書の本文を参照して、キーワード間に付与された位置関係等の検索条件を満たすものを抽出する本文サーチ手段とを備えたことを特徴とする文書検索装置。

【請求項2】文書情報を文字コードデータとして蓄積した文書データベースを対象として、検索者が指定したキーワードを含む文書を検索する文書検索装置において、該文書データベースに文書を登録する際、該登録文書の本文文字列をひらがな、漢字、及び英数字等の文字種ごとに分割し、分割した各部分文字列の間で相互に文字列の包含関係を調べ、他の文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成する手段と、該凝縮本文中に現れる文字を重複なく登録した文字成分表を作成する手段と、登録対象文書の本文のほかに凝縮本文と文字成分表を合わせて文書データベースへ格納する手段と、検索時に検索者が指定したキーワードを構成する全ての種類の文字を含む文書を、文字成分表を参照して抽出する文字成分表サーチ手段と、該文字成分表サーチで抽出された文書の件数を検査する手段と、該検査の結果件数が一定数以上の件数に達した場合は、凝縮本文を全件読み出して、検索者が指定したキーワードを構成する部分文字列を含む凝縮本文に対応する文書を抽出する凝縮本文サーチ手段と、該凝縮本文サーチ手段により抽出された文書の本文を参照して、キーワード間に付与された位置関係等の検索条件を満たす文書を抽出し、前記の文字成分表サーチの結果件数が一定数以下の場合に該文字成分表サーチで抽出された文書に対する本

文を参照して、指定キーワードを含むとともにキーワード間に付与された位置関係等の検索条件を満たす文書を抽出する本文サーチ手段とを備えたことを特徴とする文書検索装置。

【請求項3】文書情報を文字コードデータとして蓄積した文書データベースを対象として、検索者が指定したキーワードを含む文書を検索する文書検索方法において、該文書データベースに文書を登録する際、該登録文書の本文文字列をひらがな、漢字、及び英数字等の文字種ごとに分割し、分割した各部分文字列の間で相互に文字列の包含関係を調べ、他の文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成するステップと、該凝縮本文中に現れる文字を重複なく登録した文字成分表を作成するステップと、登録対象文書の本文のほかに凝縮本文と文字成分表を合わせて文書データベースへ格納するステップと、検索時に検索者が指定したキーワードを構成する全ての種類の文字を含む文書を、文字成分表を参照して抽出する文字成分表サーチステップと、該文字成分表サーチで抽出された文書の件数を検査するステップと、該検査の結果件数が所定数以上の件数に達した場合は、凝縮本文を全件読み出して、検索者が指定したキーワードを構成する部分文字列を含む凝縮本文に対応する文書を抽出し、前記検査の結果件数が所定数以下の場合は、前記文字成分表サーチで抽出された文書の凝縮本文を参照して、検索者が指定したキーワードを構成する部分文字列を含む凝縮本文に対応する文書を抽出する凝縮本文サーチステップと、抽出された文書の本文を参照して、キーワード間に付与された位置関係等の検索条件を満たすものを抽出する本文サーチステップとを有することを特徴とする文書検索方法。

【請求項4】文書情報を文字コードデータとして蓄積した文書データベースを対象として、検索者が指定したキーワードを含む文書を検索する文書検索方法において、該文書データベースに文書を登録する際、該登録文書の本文文字列をひらがな、漢字、及び英数字等の文字種ごとに分割し、分割した各部分文字列の間で相互に文字列の包含関係を調べ、他の文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成するステップと、該凝縮本文中に現れる文字を重複なく登録した文字成分表を作成するステップと、登録対象文書の本文のほかに凝縮本文と文字成分表を合わせて文書データベースへ格納するステップと、検索時に検索者が指定したキーワードを構成する全ての種類の文字を含む文書を、文字成分表を参照して抽出する文字成分表サーチステップと、該文字成分表サーチで抽出された文書の件数を検査するステップと、該検査の結果件数が一定数以上の件数に達した場合は、凝縮本文を全件読み出して、検索者が指定したキーワードを構成する部分文字列を含む凝縮本文に対応する文書を抽出する凝縮本文サーチステップと、該凝縮本文サーチステップにより抽出された文書の

本文を参照して、キーワード間に付与された位置関係等の検索条件を満たす文書を抽出し、前記の文字成分表サーチの結果件数が一定数以下の場合に該文字成分表サーチで抽出された文書に対する本文を参照して、指定キーワードを含むとともにキーワード間に付与された位置関係等の検索条件を満たす文書を抽出する本文サーチステップとを有することを特徴とする文書検索方法。

#### 【発明の詳細な説明】

##### 【0001】

【発明の属する技術分野】本発明は、文書データベースを文字列を指定して文書の全文を対象として探索する検索技術に係わり、特に補助的なファイルを用いて高速な検索処理を実現するための検索技術に関する。

##### 【0002】

【従来の技術】従来の文書検索システムでは、登録する文書の内容を表す単語（キーワードと呼ぶ）をインデクスとする方式がとられている。しかし、この方式ではインデクサーとよばれるキーワード付けの専門家が文書を逐一読み、内容を理解した上で適切なキーワードを振る必要があった。この登録時の手間の掛かる作業を回避するために、「特開昭63-198124」のような本文中に出現する単語を全てキーワードとしてインデクスファイルに登録する方法も提案されている。しかし、上記の方法ではインデクスファイルの作成時に、意味を持つ最小の単位の単語を決定するのが難しく、単語辞書あるいは、文法規則の不備のために、文章の解析に失敗して、重要な単語がキーワードとして抽出されないという問題がある。この問題を解決するために検索時に文書を文字コード化したテキストとして直接計算機に登録し、検索時にはテキストデータベース内の全ての文書の内容を読んで、与えられたキーワード（従来システムにおける統制キーワードと区別するために、以後検索タームと呼ぶ）を含む文書を探し出だすフルテキストサーチが提案されている。

【0003】このフルテキストサーチ方式は、「情報処理学会研究報告 vol. 89, no. 66 情報学基礎14-7 テキストデータベース管理システムSIGMAとその応用（1989. 7. 27）」の第2節冒頭で述べられているように、テキストファイル全体を先頭から一文字ずつ走査することが大きな特徴である。こうすることにより、キーワードに対応する文書識別子等を記述したインデクスファイルがなくとも、テキストデータベースのテキスト本体を手掛かりに検索することが可能となる。すなわち、与えられた検索タームでテキストデータ全体を文字列探索し、検索タームが記述されている文書を検索結果として出力することができる。

##### 【0004】

【発明が解決しようとする課題】従来のフルテキストサーチ方式は、テキストファイル全体を先頭から一文字ずつ走査するために処理時間が掛かり、大規模なデータベ

ースに適用すると時間が掛かるという課題があった。同文献第2節中にみられるように、汎用の大型計算機を持ってしても、2MB/s程度の検索処理速度しか実現できない。この速度でも、数メガバイト程度のデータベースであれば、検索時間は実用域内に入る。しかし、オフィス等の実用規模のデータベースには数百メガバイトの容量が必要とされ、この場合には十分な検索レスポンスが得られないことになる。

【0005】本発明の目的は、高速な文書検索方法および装置を提供することにある。

##### 【0006】

【課題を解決するための手段】上記課題を改善するために、以下の処理ステップから構成される文書検索方法を用い、該方法を実施する装置を構成する。

(1) 本文自体を格納するステップ

(2) 格納した本文を単語レベルで部分文字列へ分解し、分解した部分文字列間で相互に文字列の包含関係を調べ、他の部分文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成するステップ

(3) 本文中で用いられている文字を重複なく集めた文字成分表を作成するステップ

(4) 与えられた検索タームを文字レベルで分解し、検索タームを構成する全ての文字を含む文書を抽出する文字成分表サーチのステップ

(5) 文字成分表で抽出された文書に対応する凝縮本文を参照し、与えられた検索タームを含む文書を抽出する凝縮本文サーチのステップ

(6) 与えられた検索条件式が複数の検索ターム間の本文中での位置関係を指定している場合には、凝縮本文で抽出された文書に対応する本文データを参照し、与えられた検索タームを含み、なおかつ検索ターム間に付与された位置関係等の検索条件を満たすものを抽出する本文サーチのステップ

【作用】このように、文字成分表サーチ、凝縮本文サーチと階層的に絞り込みを行い最後に本文サーチを行う階層型プリサーチ手段を設けることによって、文字成分表サーチ、凝縮本文サーチで与えられた条件式を満たさない文書をテキスト本文を参照する以前に切り捨てて、検索対象のテキスト本文を探索する量を少なくすることができる。すなわち、検索処理時間に占める割合が高い本文検索処理時間を減らすことによって、全体の検索処理時間を短縮することが可能となる。例えば、「本文中に“画像”と“処理”とが同一の文（センテンス）内にある文書を探せ」という二つの検索タームの本文中での位置的な関係まで指定した条件式が与えられた場合、直接本文を参照する従来の方法では検索処理速度を2MB/sと仮定して、500MBのフルテキストを全て探索するのに250秒、すなわち約4分掛かる。しかし、階層型プリサーチでは、典型的な場合で、文字成分表でデータベース全件の10%に、凝縮本文でさらにその10%

に絞り込めたとすると、凝縮本文の容量が本文の30%の場合、文字成分表の容量はデータベース全体からみると無視できるほど小さいので、検索すべき凝縮本文の容量は15MBで、本文データの探索量は全データベース量の1%、すなわち5MBとなるため、2MB/sの検索速度でも、10秒で検索処理を終了できることになる。このように、「階層型プリサーチ方式」では、「文字成分表」と「凝縮本文」という2段階のプリサーチを事前に行い、それぞれ「文字レベル」と「単語レベル」のふりに掛け、最も時間を要する本文サーチの対象となる文書数をあらかじめ最小に絞り込んでおくことによって、探索文書容量を削減することができるため、等価的に非常に高速なフルテキストサーチが実現できることになる。また、条件式が単一の検索タームあるいは複数の検索タームでのAND、OR、NOT条件の場合には、凝縮本文サーチでの結果をそのまま最終検索結果とすることができる。なぜなら、凝縮本文中に存在している単語は、必ず本文中にも存在するためもう一度本文を検索する必要がないためである。このように、「単語レベル」での検索では処理時間のかかる本文サーチをまったく省略することができるため、より一層全体の検索処理時間を短縮することが可能となる。以上の処理ステップから構成されるフルテキストサーチ方法によれば、直接本文を探索する量を予め少なくすることができるため、高速なフルテキストサーチが可能となる。

#### 【0007】

【発明の実施の形態】以下、本発明の第一の実施例について図1を用いて説明する。本装置は、ディスプレイ100、キーボード101、中央制御装置CPU102、文字成分表105、凝縮本文104、及び本文103格納用ファイル110、フロッピディスクドライバ106、主メモリ200から構成される。

【0008】主メモリ200には、本文登録プログラム201、凝縮本文作成登録プログラム202、文字成分表作成登録プログラム203、文字成分表サーチプログラム204、凝縮本文サーチプログラム205、本文サーチプログラム206、階層型プリサーチ制御プログラム207が格納されるとともに、データエリア208が確保されている。これらのプログラムはCPU102で実行される。文書の登録の際は、キーボード101から入力されるコマンドにより、CPU102がフロッピディスクドライバ106に挿入されるフロッピディスクから文書データを読み込み、本文登録プログラム201を実行して読み込んだ文書データを本文103としてファイル110へ格納する。次にCPU102は、凝縮本文作成登録プログラム202を実行して、本文103を単語レベルで部分文字列へ分割し、分割した部分文字列間で相互に文字列の包含関係を調べ、他の部分文字列に含まれる文字列を排除した部分文字列の集合からなる凝縮本文を作成し、これを凝縮本文104としてファイル11

0へ格納する。最後にCPU102は、文字成分表作成登録プログラム203を実行して、本文103中で用いられている文字を重複なく集めた文字成分表を作成し、これを文字成分表105としてファイル110へ格納する。検索の際には、キーボード101から入力された検索条件式が、CPU102に送られる。CPU102では、まず階層検索制御プログラム207を実行し、その制御に基づいて文字成分表サーチプログラム204、凝縮本文サーチプログラム205、本文サーチプログラム206を順次実行する。すなわち、文字成分表サーチでは、入力された検索条件式中の検索タームを文字レベルで分解し、検索タームを構成する全ての文字を含む文書を抽出する。そして、文字成分表で抽出された文書に対応する凝縮本文を参照し、与えられた検索タームを含む文書を抽出する。もし、与えられた検索条件式中に単一の検索タームか、もしくは複数の検索ターム間の論理的な関係が指定されているのみで、本文中での位置関係までは指定されていない場合には、ここで検索を終了し、凝縮本文サーチの結果を検索結果として出力する。それ以外の場合、すなわち与えられた検索条件式中に複数の検索ターム間の本文中での位置関係が指定されている場合には、凝縮本文サーチで抽出された文書に対応する本文データを参照し、与えられた検索タームを含み、なおかつ検索ターム間に付与された位置関係等の検索条件を満たすものを抽出し、検索結果として出力する。以上が本発明の第一の実施例のフルテキストサーチ装置の概略である。

【0009】以下、本発明の特徴である文字成分表サーチ、凝縮本文サーチ、及び本文サーチと絞り込みを行う階層型プリサーチ方式の、登録及び検索方法について概略を説明する。まず「凝縮本文」と「文字成分表」の作成は、文書の登録時に自動的に行う。この処理内容を、図2に示す。本図で、登録すべき文書が入力されると、まずそのまま「本文」として格納する。

【0010】次に、この「本文」から「凝縮本文」を作成する。「凝縮本文」は、「本文」の中から文字種（漢字、ひらがな、カタカナ、英字等）ごとに文字列を分割し、繰り返し現れる言葉の重複を排除して作成される。本文が「あいまい検索のための検索技術・・・」という文書1の場合には、「検索」が重複語として切り捨てられ、「あいまい」と「検索技術」及び「のための」が「凝縮本文」として残ることになる。また、「本文」から「文字成分表」を作成する。ここでは、「本文」に現われる文字を1ビットの情報で表す。文書1の例では、「あ」と「い」があるのでそれぞれ「1」を、また「う」はないので「0」を設定する。「検」と「索」も同様にそれぞれ「1」を設定する。以下同様にして、文字成分表の該当文字部分に、「本文」にその文字がある場合には「1」を、存在しない場合には「0」を設定する。このようにして、文書の登録時に「凝縮本文」と

「文字成分表」を自動的に作成し、階層型プリサーチの準備をしておく。検索時には、図3に示すように、登録の逆の順序でこれらの補助ファイルを参照する。まず第1に、文字成分表サーチでは、文字成分表を参照し、検索ターム中の全ての文字に対応する文字成分表の該当文字部分に「1」が設定されているものを選びだす。第2に、凝縮本文サーチでは、文字成分表で選びだされた文書の凝縮本文を参照し、条件式に与えられた検索タームを含む文書を選びだす。最後に本文サーチでは、検索タームの本文での出現位置が条件式と適合するものを選びだす。本図の例では、検索「4C」理解すなわち、「検索」と「理解」が本文中で4文字以内に近接して現れるものを探せ」という条件式で検索した例を示している。結果として文書4の「検索」と「理解」が本文中で4文字離れている文書が抽出される。

【0011】以下、本実施例で用いる文字種分割・重複語排除型凝縮本文及び文字コード依存型文字成分表の作成方法と、これらを用いた階層型プリサーチの制御方法について具体的に説明する。まず最初に本実施例で用いる文字種分割・重複語排除型凝縮本文の作成方法について説明する。図4に示すように、まず本文テキストから文字種により文字列を分割する。この時の文字種とは、漢字、ひらがな、カタカナ、英字、数字、記号その他である。これらの単一文字種の連なりからなる文字列毎に本文の文字列を分割する。次に、分割した文字列のそれぞれについて、同一文書内にある他の部分文字列にその部分文字列がそっくり含まれてしまう場合、その文字列を重複文字列として凝縮本文の対象から排除する。例えば、「検索」という部分文字列は、同一文書内にある他の「知的検索技術」という部分文字列に完全に含まれるので、この「検索」は凝縮本文には登録しない。しかし、凝縮本文サーチでは、たとえ「検索」という文字列は凝縮本文に登録されていなくとも、「知的検索技術」の部分文字列としてヒットすることになる。

【0012】このように、部分文字列の重複登録を排除して、得られた部分文字列には、図5に示すように文書毎に文字列の間にセパレータを挿入する。本図では、セパレータとして記号「，」を用いている。図2、図3ではこのセパレータは記号「|」で表されているが、このセパレータは特に文字として表す必要はなく、文字に割り当てられていない特殊なコードを使用することもできる。

【0013】次に本実施例で用いる文字コード依存型文字成分表の作成方法について説明する。図6に示すように、文字コード依存型文字成分表は、文字コードによって存在を示すビット情報として、1を立てるビット位置を決定する文字成分表である。本図ではシフトJISコードを例に説明している。同図で(XXXX)Hは文字コードを16進表示したものである。例えば「検索」という文字列が文書1の本文中に存在することを示すの

に、文書1のビットリストの(8C9F)H、(8DF5)H番目に1を設定する。ビットリスト中のこの文字に対応するビット位置を文字成分表のエントリ番号と呼ぶことにする。例えば「検」のエントリ番号は(8C9F)H、または10進表示すれば35999となる。以上の文字成分表と凝縮本文を用いた階層型プリサーチの制御及びサーチ動作について説明する。まず検索条件式中の検索タームをそれぞれ一文字単位に分解し、文字成分表サーチを行う。ここでは与えられた検索タームを構成する文字コードに対応するビットリスト中のエントリ番号の位置がすべて1となるビットリストを持つ文書を求めることとなる。例えば、「検索」という文字列が検索タームとして与えられた場合、「検」、「索」に対応するビットリストの(8C9F)H、(8DF5)H番目のビットがすべて1である文書1、2、3、4、...を文字成分表サーチの検索結果とする。

【0014】すなわち、図7に示すように「検」を示す(8C9F)Hのエントリ番号のビットリスト701と、「索」を示す(8DF5)Hのエントリ番号のビットリスト702との間でビット毎にAND演算を施し、ビットAND演算結果703を得る。このビットAND演算結果703のビットリスト中で、1となっているビット位置に対応する文書番号が文字成分表サーチの検索結果としてのヒット文書を表すことになる。すなわち、「検」と「索」を全て含む文書が抽出されることになる。また、「湖」のように検索タームがただ1個の文字から構成される場合は、ここで文字成分表サーチの結果を出力して検索を終了することができる。次に文字成分表サーチで抽出された文書の凝縮本文に対してサーチを行う。ここでは図5のように文書毎に登録された凝縮本文の内容をスキャンして、与えられた検索タームを単語として含む文書を抽出する。つまり、「検」と「索」の2文字が「検索」と連続して現れる文書を抽出する。すなわち、「検」と「索」が含まれていても、「検出」と「探索」というように、別の単語として現われるようなものはここで切り捨ててしまう。このためには、文字成分表サーチで絞り込まれた文書毎の凝縮本文について本文テキストデータと同じように、一文字ずつスキャンしながら探索する。この時、文字成分表サーチで得られた結果の文書番号に対応する凝縮本文しかスキャンしない。例えば、文字成分表サーチの結果が文書番号1、2、3、4、...であれば、凝縮本文サーチでは、文書番号1、2、3、4、...の凝縮本文をスキャンする。そして、実際に凝縮本文中に検索タームが存在する文書を凝縮本文サーチの検索結果として出力する。このように、「階層型プリサーチ方式」では、「文字成分表」と「凝縮本文」という2段階のプリサーチを事前に行い、それぞれ「文字レベル」と「単語レベル」のふりに掛け、最も時間を要する本文サーチの対象となる文書数をあらかじめ最小に絞り込んでおくことによって、



探索文書容量を削減することができるため、等価的に非常に高速なフルテキストサーチが実現できることとなる。すなわち、文字成分表サーチでは、文字成分表が文字の存在を1ビットの情報で表しているため、サーチするデータ容量を極めて小さくすることができ、その結果検索時間も短時間に納めることが可能となる。さらに、キーワードを構成する文字毎のビットリストの論理積を取ることによって、キーワードに関連のない文書を大幅に切り捨て、以降の対象文書を格段に絞り込むことが可能となる。また、凝縮本文サーチでは、本文を直接スキャンするよりもデータ量が少ない分、検索処理時間が短縮できることになる。

【0015】次に、本発明の第二の実施例を説明する。本実施例は、複数の検索タームが指定された場合でも、効率的に階層型プリサーチを行うことのできるフルテキストサーチ方法を提供するものである。例えば、「“検索” AND “理解”」という条件式が与えられたときには、まず、第1ステップとして文字成分表をサーチする。ここでは与えられた検索ターム毎にそのすべての文字を含む文書を探し、その後検索ターム間で与えられた条件を満たすような文書を出力する。「“検索” AND “理解”」という条件式の場合には、“検索”の2文字を含み、かつ“理解”の2文字を含む文書を探す。すなわち、「(‘検’ AND ‘索’) AND (‘理’ AND ‘解’)」従って、「‘検’ AND ‘索’ AND ‘理’ AND ‘解’」つまり、上記の4文字を同時に含む文書を検索する。次に、この文字成分表サーチの結果絞り込まれた文書の凝縮本文をサーチする。ここでは、指定されたキーワードが単語として現われる文書だけを抽出する。すなわち、“検索”と“理解”を両方同時に含む文書を検索する。この例の場合のように、検索ターム間の関係が“AND”、“OR”等の論理条件だけで、その他にキーワード間の位置関係を規定する条件が指定されていない場合には、ここで検索を終了し、凝縮本文サーチの結果を最終検索結果として出力する。もし、位置条件が指定されている場合には、凝縮本文サーチで抽出された文書の本文をサーチし、指定条件に合致するものを抽出し、これを最終検索条件として出力する。以上が本実施例における検索動作の説明である。このように、文字成分表サーチ、凝縮本文サーチで検索ターム間の論理積を取ることで、複数の検索タームが指定された場合でも、効率的に階層型プリサーチを行い、高速なフルテキストサーチを実現することができる。

【0016】これより第三の実施例として、さらに一般的に階層型プリサーチの検索制御について詳細に説明する。図8にこのときの階層型プリサーチの制御の手順をPAD図にて説明する。ここでは「“計算機”と“知的インタフェース”のどちらかを含む文書を探せ」すなわち「“計算機” OR “知的インタフェース”」という検索式を例にあげて説明する。まず、最初にステップ80

00で文字成分表サーチを行う。ここでは与えられた検索ターム毎にそのすべての文字を含む文書を探し、その後検索ターム間で与えられた複合条件を満たすような文書を出力する。この例では、図9に示すように“計算機”を構成する3個の文字のそれぞれについて文字成分表の該当するエントリ番号間のビットAND演算を行い、次に同様に“知的インタフェース”を構成する9個の文字のそれぞれについて文字成分表の該当するエントリ番号間のビットAND演算を行い、最後に先に作成した“計算機”に対するのときのビットAND演算結果とそのビット列のOR演算を行う。すなわち、

「(‘計’ AND ‘算’ AND ‘機’) OR (‘知’ AND ‘的’ AND ‘イ’ AND ‘ン’ AND ‘タ’ AND ‘フ’ AND ‘ェ’ AND ‘ー’ AND ‘ス’)」という検索式を実行することになる。これにより、“計算機”を構成する3個の文字をすべて含む文書、もしくは“知的インタフェース”を構成する9個の文字をすべて含む文書が抽出される。以上の文字成分表サーチの結果件数が0件であれば、第8図に示すようにここで0件という検索結果を出力して検索を終了する。また、‘湖’のように検索タームがただ1個の文字から構成される場合も、ここで文字成分表サーチの結果を出力して検索を終了する。もし、検索タームが複数の文字で構成されていて、かつ文字成分表サーチの結果件数が0件でなければ、次に凝縮本文サーチを行う。凝縮本文に登録されている内容は、文字種ごとに分割された文字列である。例えば、“知的インタフェース”のように、途中で文字種が異なれば凝縮本文では部分文字列へ分解され、「知的、インタフェース」のように分割点にセパレータが入る。したがって、“知的インタフェース”のように異なる文字種から構成される検索タームの場合、このままでは凝縮本文をサーチしても該当する文字列が存在しないことになる。そこで、凝縮本文サーチに入る前に検索タームをチェックし、異なる文字種で構成される検索タームはこれを文字種毎に分割する。このように文字種で分割するという処理を施した検索タームを元々の検索タームと区別して、分割検索タームと呼ぶ。そして凝縮本文サーチは、例えば“計算機”、“知的”、“インタフェース”のように分割検索タームで検索する。ただし、分割検索タームに関しては、分割元を同じくするターム間でAND条件で検索を行う。例えば、「“計算機” OR “知的インタフェース”」という条件式の場合、凝縮本文サーチでは「“計算機” OR (“知的” AND “インタフェース”)」すなわち、「“知的”と“インタフェース”が同一文書内に存在するか、または“計算機”が存在する文書を探せ」という条件式として検索を行うことになる。凝縮本文サーチの結果が0件であれば、ここで0件という検索結果を出力して検索を終了する。また近傍条件、または文脈条件の指定の有る場合、あるいは“知的インタフェース”の

ような分割される検索タームがある場合、つまり検索タームと分割検索タームが異なる場合に限り本文サーチを行う。そうでない場合、ここで階層型プリサーチを終了し凝縮本文の結果を検索結果として出力する。ここで、文脈条件とは例えば、「“計算機” [S] “知的インタフェース”」のように示される条件式でこれは、「“計算機”と“知的インタフェース”が同一の文（センテンス）内にあるものを探せ」という意味を表す。あるいは近傍条件とは、例えば、「“計算機” [10C]

“知的インタフェース”」のように記述されるもので、これは、「“計算機”と“知的インタフェース”が10文字以内に近接して現れる文書を探せ」という意味を表す。すなわち、文脈条件、近傍条件とも文書中に出現する検索タームの位置関係を指定する検索条件のことである。このような本文中に現れる検索タームの位置関係を指定した検索条件が与えられた場合、もしくは凝縮本文中ではセパレータで区切られた途中で文字種の変わる検索タームが与えられた場合には、凝縮本文サーチの結果に対応する本文データを参照し、与えられた条件通りに本文中に検索タームが出現するものを検索結果として出力し、検索を終了することになる。このように、検索タームが異なる文字種で構成されている場合、或いは検索ターム間の本文中での出現位置に関する条件指定がある場合についても、効率的に階層型プリサーチを行い、高速なフルテキストサーチを実現することができる。

【0017】次に、本発明の第四の実施例について説明する。本実施例は、第一の実施例における文字成分表の容量を削減し、コンパクトにしたものである。第一の実施例で用いた文字コード依存型文字成分表は、処理が簡単であるが、文字成分表の1文書あたりのビットリストが長い文字成分表が大きくなるという問題がある。また、該当する文字コードが存在しないのにエントリ番号を割当てているためむだな部分が多いという問題がある。例えばシフトJISの場合、(0000)Hから(8140)Hの間、及び(A000)Hから(E040)Hの間、つまり0番目から33087番目までと40960番目から57408番目までのエントリ番号には該当する文字コードがない。それにもかかわらず、文字コードによってエントリ番号を決定するためにこの部分も全て表のエントリとして持っている必要がある。このビットリスト中のむだな部分を排除するために一旦文字コードを変換し、ビット位置を0番目からすきまなく使用できるように文字成分表を作成する。この文字コード変換型文字成分表を用いた実施例の詳細について以下説明する。

【0018】文字コード変換型文字成分表を作成するための文字コードへの変換式の例として次式をあげる。また、対応するPAD図を図10に示す。

$$\text{if } \text{SJIS} < (\text{A000})\text{H} \quad \text{then} \\ \text{SCODE} = \text{SJIS} - (\text{8040})\text{H}$$

$$\text{else} \quad \text{SCODE} = \text{SJIS} - (\text{C040})\text{H} \\ \text{SCODE} = \text{SCODE} - (\text{SCODE} / 256) \times 64 \\ \dots \dots \dots (4-1)$$

式（但し、通常文字コードの小さい値の部分は制御コードとして用いることが多いために、本式では(8140)Hとはせずに(8040)Hとして多少の余裕を持たせている。また、(SCODE/256)の演算結果の小数点以下は切り捨て、切り捨てた結果と64との乗算を行う。）式中でSJISがもとのシフトJISコードを示し、SCODEは変換後の文字コードを示す。KEISコードや他のコード体系についてもシフトJISコードとの対応がとれているので同様の式でSCODEへの変換が可能である。(4-1)式は、文字コード表に表すと図11のような変換を意味している。すなわち、(0000)Hから(FFFF)Hまでの間に(8140)H～(9FFC)H及び(E040)H～(FEFC)Hと分散して配置されている文字コードを、(0000)Hからすきまなく配置するように文字コードを変換することになる。この(4-1)式を用いてコード変換することにより、図12に示すようにビットリストの長さを非常に短くすることができ、文字成分表の全体の容量を小さくすることができる。

【0019】階層型プリサーチの制御は、第一の実施例と同じである。すなわち、図8の制御手順をそのまま使用し、第1に検索ターム中の文字を使い文字成分表サーチを行い、第2に検索タームを用いて凝縮本文サーチを行う。文脈条件の指定がなければここで検索結果を出力し、検索を終了する。文脈条件の指定があれば第3に本文サーチを行いその結果を出力する。但し、文字成分表サーチのときには入力された検索タームは全て(4-1)式に基づいて文字コード変換を施して用いることになる。以上、文字コード変換型文字成分表を用いた第四の実施例について説明した。本実施例によれば、文字コードをコード変換し、ビット位置を0番目からすきまなく並べた文字成分表を作成することにより、文字成分表の文字の割り振られていないエントリを無くすことができ、文字成分表の容量を非常に小さくすることができる。

【0020】次に、本発明の第五の実施例について説明する。本実施例は、第四の実施例における文字成分表の容量をハッシング手法を用いてさらに削減したものである。第四の実施例の文字成分表の容量をさらに小さくするために、本実施例ではビットリスト中の一つのエントリ番号に複数の文字を割り当てる。すなわち、ハッシュ関数を用いて検索ターム中の文字とビットリスト中のビット位置を対応付ける方法をとる。このハッシュ関数として例えば次の式を用いることができる。

$$\text{h}(\text{SCODE}) = \text{mod}(\text{SCODE}, \text{N}) \\ \dots \dots \dots (5-1)$$

式式中でSCODEは(4-



1) 式によってシフト JIS から変換した文字コードである。mod は第 1 引き数を第 2 引き数で割った余りを出力する関数である。N は任意の整数値である。N として、例えば 512 を用いると、‘あ’ はエントリ番号 480、‘ま’ はエントリ番号 193 となる。

【0021】このようにして作成した文字成分表の例を図 13 に示す。この場合は、N を 512 と設定したが、1 文書を登録するのに 512 ビットしか必要としないことが分かる。検索時には、与えられた検索タームの各文字について登録時と同じように、(5-1) 式のハッシュ関数を用いてエントリ番号を求め、これに対応する文字成分表のビット位置を参照する。例えば、“あいまい” という文字列の場合図 13 のようにエントリ番号 480、482、193 の位置のビットがすべて 1 の文書を文字成分表サーチの検索結果とする。こうして文字成分表サーチで求められた文書について、次にその凝縮本文をサーチする。

【0022】以下、凝縮本文サーチ及び本文サーチの制御手順について、図 14 を用いて説明する。第一の実施例では、文字成分表サーチの後検索タームが一文字からなる場合には、文字成分表サーチの結果を検索結果として出力して階層型プリサーチを終了していた。しかし、この本実施例で用いた文字成分表の文字成分表サーチでは、検索ノイズの生じる可能性があるために、凝縮本文サーチまで階層型プリサーチを継続する必要がある。例えば、ひらがなの‘は’ (シフト JIS コード (82CD) H) は、(5-1) 式でエントリ番号 13 であるが、漢字の‘艦’ (シフト JIS コード (8ACD) H) も同じエントリ番号 13 となる。このことは、検索タームとして‘艦’が与えられた場合、‘は’を含む文書もすべて文字成分表サーチの結果として検索されてくることになる。したがってさらに、凝縮本文をスキャンして実際に漢字の‘艦’を含む文書を抽出し、これを検索結果として出力することになる。以上、第五の実施例について説明した。本実施例ではハッシュ関数を使って、文字成分表の 1 エントリに複数の文字を割り当てることにより、文字成分表の容量を格段に小さくできるという効果が得られる。

【0023】次に第六の実施例について説明する。第五の実施例のように単純にハッシングした場合、ひらがなのように文書中に出現しやすい文字と、JIS 第 2 水準の漢字のようにめったに出現しない文字とが同じエントリ番号となる可能性がでてくる。例えば、ひらがなの‘は’と、漢字の‘艦’は同じエントリ番号 13 となり、検索タームとして‘艦’が与えられた場合‘は’を含む文書はすべて文字成分表サーチの結果としてヒットすることになる。ひらがなの‘は’は日本語の文書では非常に使用頻度の高い文字のためほぼ全件の文書が文字成分表サーチでヒットする。このように文字成分表サーチでの絞り込みの率が低下すると、凝縮本文もスキャン

する文書量が増えるために全体の検索処理時間が増大することになる。

【0024】このような絞り込み率の低下を防ぐためには、ハッシュ関数を文字の使用頻度を考慮して定める必要がある。本実施例で用いる文字成分表を文字種別ハッシング型文字成分表と呼ぶ。文字種別ハッシング型文字成分表を作成するには、例えば図 15 に示すように、各文字種毎に文字成分表のエントリ領域を割り当て、その領域内で文字コードにより折り返すようなハッシュ関数を作る。このようなハッシュ関数を実現するには、文字コードによって文字種を判定した後、mod 関数で折り返してもよいし、文字コードとエントリ番号との対応表により実現することもできる。このハッシュ関数の一例を図 16 に PAD 図で示す。本実施例では、ひらがな、カタカナ、英字のエントリ数をそれぞれ 20 とし、記号のエントリ数を 10、数字のエントリ数を 10、JIS 第 1 水準のエントリ数を 370、JIS 第 2 水準のエントリ数を 61 としている。まず、入力された検索タームに対して、文字コードにより文字種を判定し、それぞれの文字種ごとに文字成分表の割り当てられたエントリの部分を mod 関数を用いて折り返す。

【0025】すなわち、SCODE が (01DF) H から (0231) H の範囲にあれば、ひらがな文字列であるので、その SCODE を 20 で mod をとってこれをエントリ番号とする。SCODE が (0240) H から (0296) H の範囲にあれば、カタカナ文字列であるので、その SCODE を 20 で mod をとって、これにカタカナのハッシング領域の先頭である 20 を足した値をエントリ番号とする。SCODE が (01A0) H から (01DA) H の範囲にあれば、英字文字列であるので、その SCODE を 20 で mod をとって、これに英字のハッシング領域の先頭である 40 を足した値をエントリ番号とする。SCODE が (018F) H から (0198) H の範囲にあれば、数字文字列であるので、その SCODE を 10 で mod をとって、これに数字のハッシング領域の先頭である 70 を足した値をエントリ番号とする。SCODE が (065F) H から (1232) H の範囲にあれば、JIS 第 1 水準の漢字文字列であるので、その SCODE を 370 で mod をとって、これに JIS 第 1 水準の漢字文字列のハッシング領域の先頭である 80 を足した値をエントリ番号とする。SCODE が (125F) H から (1FDE) H の範囲にあれば、JIS 第 2 水準の漢字文字列であるので、その SCODE を 61 で mod をとって、これに JIS 第 2 水準の漢字文字列のハッシング領域の先頭である 450 を足した値をエントリ番号とする。以上の SCODE 以外の場合には、記号その他の文字種による文字列とみなし、その SCODE を 10 で mod をとって、これに記号のハッシング領域の先頭である 60 を足した値をエントリ番号とする。

【0026】この文字種別ハッシング型文字成分表を用いた階層型プリサーチの制御手順は、第五の実施例と同じである。すなわち、第1に検索ターム中の文字を用いて文字成分表サーチを行い、第2に検索タームを用いて凝縮本文サーチを行う。文脈条件等が指定されていない場合には、ここで検索を終了するが、そうでない場合には、第3に本文サーチを行い結果を出力する。以上説明したように、本実施例によれば、使用頻度を考慮して文字種ごとに文字成分表のエントリ番号を対応させた文字種別ハッシング型文字成分表を用いることにより、文字成分表サーチでのノイズを少なくできるため、凝縮本文における文書のスキャン量が減り、その分高速なフルテキストサーチが可能となる。

【0027】次に第七の実施例として、さらに文字成分表サーチにおける絞り込みの率を向上させ、凝縮本文のスキャン量を減らすことのできる頻度情報ハッシング型文字成分表を用いた階層型プリサーチの制御方法を説明する。頻度情報ハッシング型文字成分表を作成するには、データベースに登録してある文書の文字の使用頻度を調べ、頻度情報によりハッシュ関数を決定する。頻度の大きい文字については、同一エントリにできるだけ他の文字が入らないようにし、頻度の少ない文字について同一エントリに複数の文字が入るようにハッシュ関数を調整する。こうすることにより、平均的に安定した絞り込み率が文字成分表サーチで得られることになる。具体的には、図17に示すように(4-1)式で得られるSCODEをもとに一度データベース中で該当する文字を使用している文書数を調べ頻度順に並べ替える。次に、頻度の大きいものから文字成分表のエントリ数分 $N_t$ だけとる。そして $N_t$ 以内の頻度数分布のうち最も上位の頻度を持つエントリだけを残して、その他のエントリに順次 $N_t$ 以上のエントリ番号を割り付けていく。この $N_t$ 以上のエントリ番号の割り付けには( $N_t+1$ )番目のエントリを $N_t$ のエントリとし、( $N_t+2$ )番目を( $N_t-1$ )番目のエントリとするように、 $N_t$ より順次頻度の大きいエントリを割り付けていく。割り付けていく過程では、常に最上位の頻度を持つエントリの上には、他のエントリを割り付けないようにする。割り付けたエントリは、図18に示すようにテーブルの形で、記憶しておきこのテーブルを参照してハッシュ関数を構成する。すなわち、SCODEが(095F)Hの文字‘検’は、エントリ番号231であることが分かる。

【0028】階層型プリサーチの制御手順は、第五の実施例と同じである。すなわち、図14の制御手順をそのまま使用し、第1に検索ターム中の文字を用いて文字成分表サーチを行い、第2に検索タームを用いて凝縮本文サーチを行う。文脈条件等が指定されていない場合には、ここで検索を終了するが、そうでない場合には、第3に本文サーチを行い結果を出力する。このように、本実施例によれば、データベース中で実際に用いられる文

字の頻度分布をもとに文字成分表を作成することによって、文字成分表サーチで常に安定して高い絞り込み率が得られるため、検索タームに依存せず安定して短時間の検索処理時間を得ることができる。

【0029】以上、文字成分表の異なる実施例について五つの実施例を説明した。これより凝縮本文の異なる実施例についての説明をする。第一の実施例で用いた凝縮本文は作成の処理が簡単であるが、図4でも分かるように“のための”というような本来検索に使われないような文字列まで凝縮本文に残ることになる。このことは凝縮本文の圧縮率低下を招く。つまり、検索時にスキャンする凝縮本文の量が増えるため、検索処理時間が増加してしまう。このような、凝縮本文の圧縮率を低下させる主な要因は、“のための”というような付属語の連なったそれ自体では意味を持たない文字列を凝縮本文に登録してしまうところにある。

【0030】そこで、第八の実施例として、この検索に不要な付属語の連なりを除去した凝縮本文を用いる階層型プリサーチを説明する。この凝縮本文を文字種分割・重複排除・付属語除去型凝縮本文と呼ぶ。この凝縮本文の作成方法は図19に示すように、本文のテキスト文字列から文字種分割して部分文字列に分け、それから重複語を排除した後、付属語の除去を行う。文字種分割と重複排除は第一の実施例と変わらない。付属語除去は、重複排除の済んだひらがな文字列に対して行う。この付属語除去のための解析は、図20に示すように基本単語辞書と単語間の接続規則を基に行う。基本単語辞書には、図21のようにひらがなのみから構成される動詞、指示代名詞、形容詞、形容動詞、副詞、接続詞、助詞、助動詞、またこれらの品詞の活用語尾が品詞情報とともに登録されている。本図の例では、動詞として<ある>、<なる>、<もつ>等がそれらの活用語尾とともに登録されている。接続規則には基本単語辞書に登録された各語が他のどの語と接続し得るかを登録する。例えば図22に示すように、<動詞-もつ連体形>に<名詞-こと>が接続し、さらに<名詞-こと>には<助詞-が>が接続し得ることが登録されている。このような基本単語辞書及び接続規則を用いてひらがなの部分文字列が付属語から構成されているか否かを判定し、凝縮本文へその文字列を登録するか否かを決定する。例えば、“のための”という部分文字列は<助詞-の><名詞-ため><助詞-の>というように接続した文字列と解析できるため、付属語のみから構成された文字列と判定し排除する。一方、“あいまい”という文字列は、付属語と解析ができないため排除せずにそのまま凝縮本文へ登録する。

【0031】このように、付属語を解析してひらがな文字列を排除し、検索に使われることのない無用の情報を削除することによって、凝縮本文の圧縮率を高めることが可能となる。また解析に用いる基本単語辞書と接続規

則は、時代とともに登録語数が増えていく従来のキーワード辞書とは基本的に異り、普遍的なもので一度作成してしまえば更新していく必要がないという利点がある。付属語として解析できるものだけを排除するために、辞書に存在しないひらがなから構成される新語が現れても必ず凝縮本文に残るということになる。

【0032】次に、文字種分割・重複排除・付属語除去型凝縮本文を用いた階層型プリサーチ方式の制御について説明する。文字種分割・重複排除・付属語除去型凝縮本文では、ひらがな文字列を付属語解析して凝縮本文に登録しない場合がある。そのため、特定のひらがな文字列で検索しようとした場合、凝縮本文サーチで検索もれとなる場合がある。例えば“めまい”という文字列は、動詞の未然形活用語尾“め”と助動詞“まい”の終止形が接続したものと解析できる。具体例としては、“認めまい”があげられる。ところが“めまい”は、名詞として使われている場合でも、付属語除去処理の結果凝縮本文からは削除されてしまう。したがってこのような場合、“めまい”で凝縮本文を検索すると検索もれが生じる可能性がでてくる。そのため、検索タームが凝縮本文中にもともとない言葉なのか、あるいは凝縮本文作成過程で除去された可能性のある言葉なのかをチェックしてから検索する必要が生じる。検索タームが凝縮本文に登録されるべき語か否かというチェックは、凝縮本文を作成したときに用いた付属語除去のアルゴリズムをそのまま適用する。この例では、“めまい”という検索タームが与えられたときは、これが付属語の連なりと判定することができる。

【0033】以上の検索制御の手順を図23で説明する。まず文字成分表サーチを行う。結果件数が0件であれば、0件を検索結果として出力して検索処理を終了する。第一の実施例でも述べたが、ハッシュ関数を用いない方式では検索タームが一文字の場合にかぎり、文字成分表のサーチ結果を最終検索結果として出力できる。すなわち、第一及び第四の実施例で説明した文字成分表を用いる場合には、検索タームが一文字であるか否かを調べ、一文字であれば文字成分表サーチの結果を検索結果として出力し、処理を終了する。第五、第六、第七の実施例で述べたハッシュ関数による文字成分表を用いる場合には、この検索タームが一文字か否かというチェックは行わず、常に次の凝縮本文サーチを行う。この後、第一の実施例と同様に、分割検索タームを生成する。

【0034】次に、分割検索タームのそれぞれについて付属語解析を行う。分割検索タームのうち一つでも付属語と判定された場合、その分割検索タームは凝縮本文から削除されている可能性があるため、凝縮本文サーチを行わず、文字成分表サーチの結果に基づいて本文を直接サーチする。一方、付属語解析の結果、分割検索タームが全て付属語でないと判定されたならば、第一の実施例と同様に凝縮本文サーチを行う。近傍条件あるいは、文

脈条件の指定がない場合、あるいは分割検索タームがもとの検索タームと同じ場合には、この凝縮本文サーチの結果を最終検索結果として出力し、検索を終了する。もし、近傍条件ないし文脈条件が指定されている場合、あるいは分割検索タームと元の検索タームが異なる場合には、次に本文サーチを実行し、その結果を最終的な検索結果出力とする。このように、本実施例によれば、ひらがな文字列を解析し、不要な付属語の連なりを凝縮本文から除去した文字種分割・重複排除・付属語除去型凝縮本文を用いることにより、凝縮本文の圧縮率を向上させ、検索処理時間を短縮することができる。

【0035】次に、第九の実施例として、ひらがな文字列を全て排除した、文字種分割・重複排除・ひらがな文字列除去型凝縮本文を用いた階層型プリサーチを説明する。第八の実施例で説明した凝縮本文は、確かに圧縮率が上がるものの付属語解析の際に誤った解析をする可能性がある。例えば第八の実施例でも用いた“めまい”という文字列の例の外にも、付属語解析だけでは本質的にどれが付属語か正しく判定できない場合がまれにある。例えば、“動作してこの応用...”という文書の中の“してこの”という部分文字列は、“～して、この～”という意味で用いられているのか、“～し、てこの～”のように機械のてこを意味しているのかが判定するのが難しい。後者の意味で用いられている場合には、“てこ”という検索タームを指定した際に、“てこ”は付属語と判定されないため、凝縮本文をサーチしにくくなることになる。一方、凝縮本文作成では、“してこの”が付属語と解析され凝縮本文から削除されているため凝縮本文サーチで検索もれとなってしまう。この付属語解析の不完全さを補正するために、ひらがな文字列か否かという単純な判定方法で階層型プリサーチを実現するのが、本第九の実施例である。この凝縮本文の作成方法を、図24に示す。本方法では文字種分割の後、ひらがなを除去して重複登録排除を行う。

【0036】この文字種分割・重複排除・ひらがな文字列除去型凝縮本文を用いた階層型プリサーチの制御手順について図25を用いて説明する。まず第八の実施例と同様に文字成分表サーチを行う。この後、分割検索タームを生成する。次に、分割検索タームのそれぞれについてひらがな文字列か否かチェックを行う。分割検索タームのうち一つでもひらがな文字列がある場合、凝縮本文サーチを行わず、文字成分表サーチの結果に基づいて本文を直接サーチする。一方、分割検索ターム中にひらがな文字列がない場合、第一の実施例と同様に凝縮本文サーチを行い、近傍、文脈条件の指定がある場合、あるいは分割検索タームが元の検索タームと異なる場合には、本文サーチまで検索処理を続行する。このように、本実施例によれば、ひらがな文字列を全て排除した凝縮本文を用いることによって、ひらがな文字列についても検索もれのない正確なフルテキストサーチが実現できる。

【0037】次に、本発明の第十の実施例について、説明する。上記第九の実施例では、ひらがなの検索タームが与えられた場合、本文を直接参照する必要がある。したがって検索時間がより多く掛かることになる。そこで、ひらがなの検索タームが与えられた場合でも高速にフルテキストサーチできる方法として、第十の実施例の説明をする。本実施例では、第九の実施例で用いた凝縮本文の外に第九の実施例では除去したひらがな文字列を登録した凝縮本文を別に作成する。図26に示すように、文字種分割、重複登録排除の後、残った部分文字列がひらがな文字列か否かを判定し、ひらがな文字列以外を凝縮本文Aとして登録し、ひらがな文字列を凝縮本文Bとして登録する。こうすれば、ひらがなだけの検索タームが与えられた際、凝縮本文Bを探索することができるようになるため、検索時間を短縮することが可能となる。実際の階層型プリサーチの検索制御の手順を図27に示す。まず第八の実施例と同様に文字成分表サーチを行う。もし、検索結果が0件なら、ここで検索を終了する。この後、分割検索タームを生成する。次に、分割検索タームをひらがな文字列のタームとそれ以外の文字列からなるタームに分類する。その後、ひらがな以外の文字列からなる分割検索タームがある場合には、凝縮Aをサーチする。次にひらがなの分割検索タームがある場合には、凝縮Bをサーチする。その後は、第一の実施例と同様に、近傍、文脈条件の指定がある場合、あるいは分割検索タームがもとの検索タームと異なる場合には、本文サーチまで検索処理を続行する。このように、ひらがなのみの凝縮本文と、ひらがな以外の凝縮本文と分けて格納することにより、どんな文字種の検索タームが入力されても、凝縮本文を有効に活用でき、常に高速なフルテキストサーチが実現できる。

【0038】次に、第十一の実施例について説明する。本実施例は、凝縮本文の圧縮率を上げるために、文字種毎に独立した凝縮本文を用いる方法に基づいたものである。本実施例で用いる凝縮本文を文字種分割・重複排除・文字種別登録型凝縮本文と呼ぶ。この文字種分割・重複排除・文字種別登録型凝縮本文を作成するには、図28に示すように、文字種分割、重複登録排除を行った後、残った部分文字列の文字種を判定してひらがな凝縮本文H、カタカナ凝縮本文I、漢字凝縮本文J、英字凝縮本文K、数字凝縮本文L、記号その他の文字種凝縮本文Mに分類して登録する。こうすることにより、例えば漢字の検索タームで検索する場合には、漢字文字種の凝縮本文Jのみをサーチすればよいことになるため、検索時間をさらに短縮することができる。具体的な階層型プリサーチの制御手順を図29を用いて説明する。まず、第八の実施例と同様に文字成分表サーチを行う。検索結果件数が0件なら、ここで検索を終了する。この後、分割検索タームを生成する。次に、分割検索タームを文字種毎に分類する。その後、ひらがなの分割検索タームが

ある場合には凝縮Hを、カタカナの分割検索タームがある場合には凝縮Iを、というように分解検索タームの文字種にしたがってサーチする凝縮本文を選択する。その後は、第一の実施例と同様に、近傍、文脈条件の指定がある場合、あるいは分割検索タームがもとの検索タームと異なる場合には、本文サーチまで検索処理を続行する。このように、文字種ごとに凝縮本文ファイルを分離し個々の凝縮本文の容量を小さくすることにより、漢字のみ、カタカナのみ、あるいはひらがなのみ、といった単一文字種の検索タームでのフルテキストサーチが高速に行えるという効果が得られる。

【0039】次に第十二の実施例について、図30および図31を用いて説明する。本実施例は、特願平02-193015で提案した文書検索装置を用い、本発明を実現したものである。本装置の主な構成は、キーボード3001、検索式解析プログラム3002、ビットサーチプロセッサ3007a、ストリングサーチエンジン3006、複合条件判定用マイクロプロセッサ3045a、検索結果格納メモリ3046、ディスプレイ3020、半導体メモリ装置3010a、RAMディスク装置3010b、集合型磁気ディスク3010c、及び検索実行制御プログラム3008よりなる。半導体メモリ装置3010aには文字成分表が、RAMディスク装置3010bには凝縮本文、集合型磁気ディスク3010cには本文がそれぞれ格納されている。但し、文字成分表及び凝縮本文は、集合型磁気ディスク3010cに格納されていて、本装置の運用開始時点でそれぞれ半導体メモリ装置3010a及びRAMディスク装置3010bへローディングされる。

【0040】階層プリサーチ制御の手順は、いままで実施例で説明してきたものと変わらない。いままでの実施例との相違点は、文字成分表を半導体メモリ、凝縮本文をRAMディスク、本文を集合型磁気ディスクに格納したところと、文字成分表サーチ専用のマイクロプロセッサ、凝縮本文サーチ及び本文サーチ専用のストリングサーチエンジンを用いていることである。検索処理の手順を以下に説明する。

【0041】キーボード3001から入力した検索条件式はサーチマシン制御用マイクロプロセッサMPU03050上の検索式解析プログラム3002により解析される。すなわち、検索式解析プログラム3002では検索条件式を構成するキーワード部分とそれらの包含条件及び配置条件を記述した複号条件記述部に分離する。包含条件は論理条件として記述され、配置条件は近傍条件や文脈条件として記述されたものである。分離抽出後、キーワード部分は同じくMPU03050上の同義語展開プログラム3003に渡され、複号条件記述部は複号条件解析プログラム3041に渡される。同義語展開プログラム3003では、ここに内蔵された同義語辞書を参照して、入力されたキーワードの同義語が求められ

る。そして、ここで同義語展開されたキーワード群は異表記展開プログラム3004へ渡される。本図の例の場合、“計算機”から、“電算機”、“コンピュータ”、“COMPUTER”などが生成される。異表記展開プログラム3004では、ここに入力されてきたキーワード群に対して異表記展開処理が施される。本図の例の場合、“コンピュータ”から“コンピューター”が、“COMPUTER”から“Computer”などが生成される。こうして同義語及び異表記展開されたキーワード群は、次にオートマトン生成用マイクロプロセッサMPU13005a上のオートマトン生成用プログラム3005に送られる。オートマトン生成用プログラム3005では、異表記展開プログラム3004から送られてきたキーワード群に対して、これらを一括照合するオートマトンを生成し、状態遷移テーブルと照合すべきキーワードの識別コード情報として、サーチエンジン3006に設定する。サーチエンジン3006は有限オートマトン方式に基づく高速多重文字照合回路である。また、異表記展開プログラム3004で異表記展開されたキーワード群は、該当キーワードと共に、ビットサーチ用マイクロプロセッサMPU33007a上のビットサーチプログラム3007へ渡される。

【0042】一方、近傍条件、文脈条件や、AND、OR等の論理条件は検索式解析プログラム3002から、複合条件解析プログラム3041、近傍条件解析プログラム3042、文脈条件解析プログラム3043、論理条件解析プログラム3044を経て複合条件判定プログラム3045へと送られる。必要な検索情報がビットサーチプログラム3007、ストリングサーチエンジン3006、複合条件判定プログラム3045へ送られた後、検索制御実行プログラム3008は、まずビットサーチプログラム3007に起動を掛ける。ビットサーチプログラム3007は、半導体メモリ装置3010aに格納してある文字成分表を読み出し、文字成分表サーチを行う。文字成分表サーチの結果は、検索結果格納メモリ3046へ格納する。

【0043】文字成分表サーチが終わった後、検索実行制御プログラム3008は、検索結果格納メモリ3046を参照し、検索結果が0件であれば、0件を検索結果として出力し検索処理を中断する。検索結果が0件でなければ、ストリングサーチエンジン3006へ起動をかけると同時に検索結果格納メモリ3046に格納されている文字成分表サーチの結果でヒットした文書の凝縮本文をRAMディスク装置2910bから読み出し、ストリングサーチエンジン3006へ送り、凝縮本文サーチを実行させる。この結果件数が0件であるか否かの条件判定は検索実行制御プログラム3008で行う。ストリングサーチエンジン3006では、RAMディスク装置3010bより読み出された、凝縮本文を分割検索タームでサーチする。照合結果は複合条件判定プログラム30

45に順次送られる。複合条件判定プログラム3045では、検索ターム間に付与された論理条件を判定し、条件に適合する文書の文書番号を検索結果格納メモリ3046へ順次格納する。

【0044】凝縮本文サーチが終了した後、検索実行制御プログラム3008は、もう一度検索結果格納メモリ3046を参照し、結果件数が0件であれば、0件を検索結果として出力し、検索を終了する。0件でない場合で、近傍、文脈条件が設定されているか、もしくは分割検索タームが検索タームと異なっている場合にかぎり検索結果格納メモリから、検索結果文書番号を読み取り、これに対応する本文を集合型磁気ディスク装置3010cから読み出し、ストリングサーチエンジン3006へ送り、今度は本文サーチを実行させる。近傍、文脈条件が設定されてなく、かつ分割検索タームが検索タームと等しい場合には、検索結果格納メモリに格納されている検索結果件数を出力し、検索を終了する。

【0045】ストリングサーチエンジン3006では、集合型磁気ディスク装置3010cから読み出された本文をスキャンして本文サーチを行う。結果は複合条件判定プログラム3045に順次送られる。複合条件判定プログラム3045では、検索ターム間に付与された論理条件のほか近傍、文脈条件を判定し、条件に適合する文書の文書番号を順次検索結果格納メモリ3046へ格納する。本文サーチまで実行した場合は、本文サーチの終了後、検索実行制御プログラム3008は、検索結果格納メモリ3046を参照し検索結果件数を出力して検索を終了する。このように、容量の大きな本文データを磁気ディスクに、容量の小さな文字成分表や凝縮本文を、半導体メモリやRAMディスクに格納することにより、大規模なデータベースに対しても高速なフルテキストサーチを実現することが可能となる。

【0046】次に凝縮本文を磁気ディスクに格納した第十三の実施例について説明する。凝縮本文を磁気ディスクに格納する場合、階層型プリサーチの制御の手順を最適化することによって、同一の構成を用いた通常の階層型プリサーチを実行するよりも高速に処理することができる。以下、この制御の手順について説明する。磁気ディスクは通常、機械的に動く磁気ヘッドを持っている。このため、ディスク上の情報を飛び飛びに読み出す（スキップアクセスと呼ぶ）よりも、まとまった情報を一括して読み出す（シーケンシャルアクセスと呼ぶ）方が速いという特徴がある。いま、スキップアクセスの読み出し速度を $V_{skip}$  MB/s、シーケンシャルアクセスの読み出し速度を $V_{seq}$  MB/sとすると、データベース全件の文書数を $N_a$ 件、文字成分表サーチの結果件数を $N_c$ 件とし、文書の容量が均一であるとした場合、 $N_c > (V_{skip}/V_{seq}) \cdot N_a$  ……(12-1)式のと看、シーケンシャルアクセスにより凝縮本文を全件サーチした方が、文字成分表サーチ



の結果に基づいてスキップアクセスするよりも処理時間が短くなる。したがって、図32に示すように文字成分表サーチの後、階層プリサーチ制御プログラムにおいて結果件数を判定し、(12-1)式を満たすヒット件数に達した場合には、文字成分表サーチの結果を無視して、凝縮本文をデータベース全件分サーチする。以上の方法を用いると、磁気ディスクに凝縮本文を格納するために、大容量のRAMディスクを使用しなくとも済み、比較的高速なフルテキストサーチを低価格の文書検索装置で実現できることになる。

$$N_c Q_{sr} / V_{sr} + N_{sr} Q_{tx} / V_{tx} > N_c Q_{tx} / V_{tx} \quad \dots\dots\dots (13-1) \text{ 式}$$

のとき、凝縮本文サーチをせずに、本文サーチを直接行ったほうが検索時間が短くなる。 $N_{sr}$  は凝縮本文を実際にサーチするまでわからないが、あらかじめ定数を

$$N_{sr} = \alpha N_a \quad (0 < \alpha < 1) \quad \dots\dots\dots (13-2) \text{ 式として、}$$

$$(13-1) \text{ 式を変形すると、} \\ N_c < \alpha N_a (Q_{tx} / V_{tx}) / (Q_{tx} / V_{tx} - Q_{sr} / V_{sr}) \quad \dots\dots\dots (13-3) \text{ 式}$$

のとき、本文サーチを直接行うことにする。 $\alpha$  をしきい値として検索前にあらかじめ値を設定しておき、文字成分表サーチの後(13-3)式により凝縮本文サーチを行うか否か決定する。この制御を行うことにより、近傍、文脈条件の指定の下で高速なフルテキストサーチを実現することができる。以上、第十二の実施例の廉価版のシステム構成でフルテキストサーチを実現する第十三、第十四の実施例について説明した。

【0048】このほかにも、凝縮本文をまったく使用せず凝縮本文サーチのステップを省いて、文字成分表サーチから直接本文サーチを実行する制御方法によっても階層型プリサーチを実現することができる。この方法によれば、本文をスキャンする量が増えるため検索時間は多少掛かるが、高価なRAMディスクを使用しなくとも済み、また凝縮本文を格納する磁気ディスク容量が不要となるため、さらに低価格の文書検索装置を実現できることになる。また、文字成分表を使用せずに直接RAMディスクあるいは磁気ディスク上の凝縮本文を全件サーチし、近傍、文脈条件などの検索ターム間の位置関係の検索条件指定があるときにのみ本文サーチする制御方法によっても階層型プリサーチを実現することができる。この方法によれば、凝縮本文の探索量が増えるため検索時間は多少掛かるが、文字成分表を格納する半導体メモリが不要となるため、その分低価格の文書検索装置を実現できることになる。

【0049】あるいは、今までの実施例で用いていたビットリスト形式の文字成分表を図33に示すように、文書中に現れる文字を書き列ねた形式、すなわち1文字を1ビットとして表すのではなく、そのまま文字コード自体として格納した文字成分表を使用することもできる。あるいはこの時に、第五の実施例、第六の実施例、及び

【0047】次に凝縮本文を磁気ディスクに格納した第十四の実施例について説明する。近傍、文脈条件が指定されている場合には、文字成分表サーチ結果が非常に少ない場合、凝縮本文サーチを行わずに、文字成分表サーチ結果をもとに本文を直接サーチするほうが検索時間が短くなる。今、凝縮本文のサーチ速度を  $V_{sr}$  MB/s、本文のサーチ速度を  $V_{tx}$  MB/s とし、文字成分表の結果件数を  $N_c$ 、凝縮本文の結果件数を  $N_{sr}$ 、凝縮本文の1件当たりのデータ容量を  $Q_{sr}$ 、本文の1件当たりのデータ容量を  $Q_{tx}$  とすると、

設定して凝縮本文サーチを行うか否か決定することになる。たとえば、データベース全体の文書数を  $N_a$  として

第七の実施例で説明したハッシュ関数を用いて一つの文字エントリに複数個の文字を対応させ文字成分表の容量を削減することもできる。このように文字コードを格納した文字成分表を用いた文字成分表サーチは、凝縮本文や本文サーチと同様に、一文字ずつファイルからデータを読み出し該当する文字が存在するか否か判定することで実現できる。このように、本文中で用いられている文字のみを集めた文字成分表を用いることにより、データ構造を簡素化でき、かつビット演算をせずに凝縮本文、本文サーチと同じスキャン型のサーチを用いることができるため、検索処理方法が簡素化できるという効果が得られる。

【0050】さらに、文字成分表も磁気ディスクに格納した構成でも、階層型プリサーチを実現することができる。この磁気ディスクに文字成分表を格納した場合には、文字成分表サーチにおいて検索ターム中で用いられている文字のビットリストを磁気ディスクから順次読み出しながらビット演算処理を行っていく。もしくは、上記の文字コードをそのまま文字成分表とした場合には、文字成分表を順次読み出しながら該当する文字を全て含む文書を選び出す。この文字成分表を磁気ディスクに格納する方法によれば、半導体メモリを使わずに済むために、さらに低価格の文書検索装置を実現することが可能となる。

#### 【0051】

【発明の効果】本発明によれば、文字成分表及び凝縮本文を用いて、階層的に文字レベル及び単語レベルで入力された検索タームに関連しない文書をふるい落とすことにより、無用の本文サーチを省くことができるため、高速な文書検索フルテキストサーチの実現技術となり、大規模な文書データベースでも実用的な応答速度で文書検索



を行うことが可能となる。

【図面の簡単な説明】

【図 1】本発明の第一の実施例の構成を示す図である。

【図 2】本発明の特徴となる階層型プリサーチのための登録処理を示す図である。

【図 3】本発明の特徴となる階層型プリサーチの検索処理を示す図である。

【図 4】凝縮本文を作成する一例を示した図である。

【図 5】凝縮本文の格納形態を示す図である。

【図 6】文字成分表の概要を示す図である。

【図 7】文字成分表サーチの概要を示す図である。

【図 8】階層型プリサーチの処理手順を示す図である。

【図 9】第三の実施例における文字成分表サーチの処理を示す図である。

【図 10】第四の実施例で用いる文字成分表のコード変換の処理を示す PAD 図である。

【図 11】第四の実施例で用いる文字成分表のコード変換の概要を示す図である。

【図 12】第四の実施例で用いる文字成分表の概要を示す図である。

【図 13】第五の実施例で用いる文字成分表の概要を示す図である。

【図 14】第五の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 15】第六の実施例で用いる文字成分表の概要を示す図である。

【図 16】第六の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 17】第七の実施例で用いる文字成分表の作成方法の概要を示す図である。

【図 18】第七の実施例で用いる文字成分表のためのハッシュ関数で用いる文字コード・エントリ番号の対応表の概要を示す図である。

【図 19】第八の実施例で用いる凝縮本文の作成する方法を示す図である。

【図 20】第八の実施例で用いる凝縮本文のためのひらがな文字列の処理方法を示す図である。

【図 21】第八の実施例で用いる付属語解析のための基本単語辞書を示す図である。

【図 22】第八の実施例で用いる付属語解析のための接続規則を示す図である。

【図 23】第八の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 24】第九の実施例で用いる凝縮本文の作成する方法を示す図である。

【図 25】第九の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 26】第十の実施例で用いる凝縮本文の作成する方法を示す図である。

【図 27】第十の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 28】第十一の実施例で用いる凝縮本文の作成する方法を示す図である。

【図 29】第十一の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 30】第十二の実施例の構成の部分を示す図である。

【図 31】第十二の実施例の構成の残りの部分を示す図である。

【図 32】第十二の実施例で用いる階層型プリサーチの処理手順を示す図である。

【図 33】文字として格納した文字成分表の概要を示す図である。

【符号の説明】

103 : 本文

104 : 凝縮本文

105 : 文字成分表

201 : 本文登録プログラム

202 : 凝縮本文作成登録プログラム

203 : 文字成分表作成登録プログラム

204 : 文字成分表作成検索プログラム

205 : 凝縮本文作成作成プログラム

206 : 本文サーチプログラム

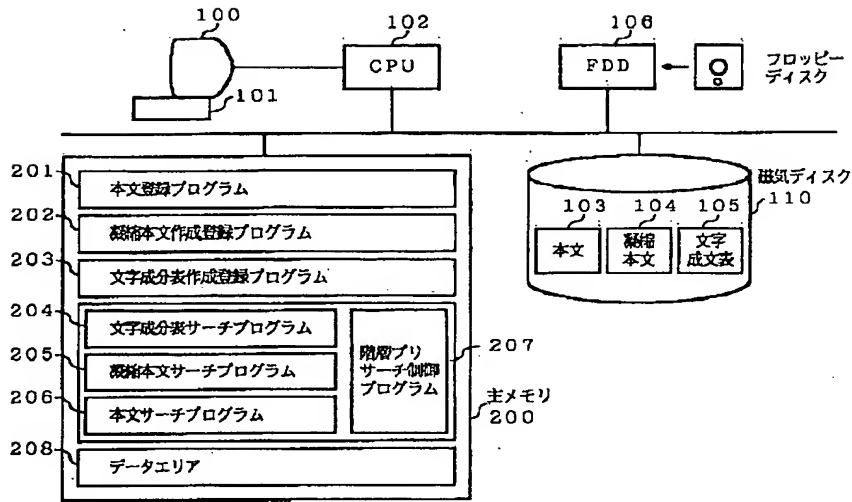
【図1】

【図5】

【図1】

【図5】

〈凝縮本文の構造〉



文書1	あいまい,のための,検索技術
文書2	自然語,による,検索技術
文書3	壁,を,検出,しながら,出口,探索
文書4	文書理解,を,用,いた,検索,システム
...	...
文書N	

【図2】

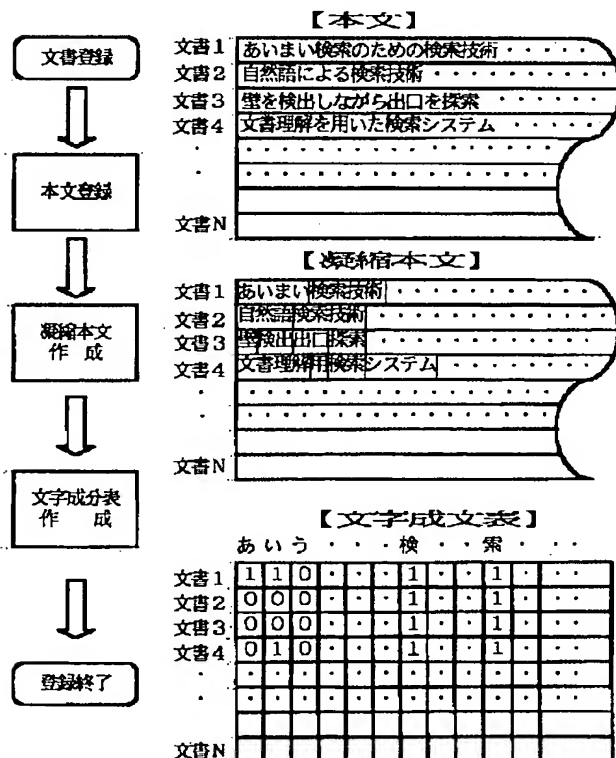
【図3】

【図2】

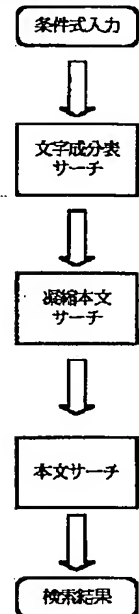
〈階層型プリサーチにおける文書の登録処理方法〉

【図3】

〈階層型プリサーチにおける検索処理方法〉



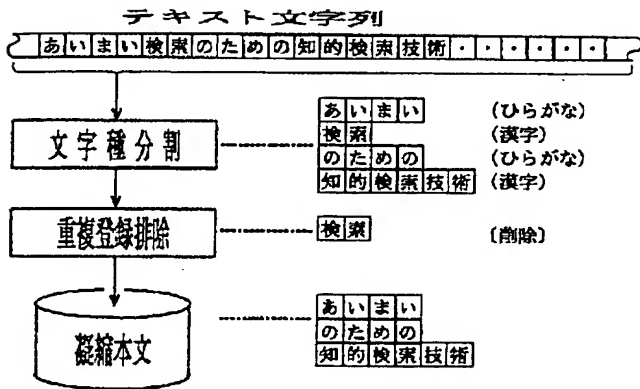
“検索 [4C] 理解”



【図4】

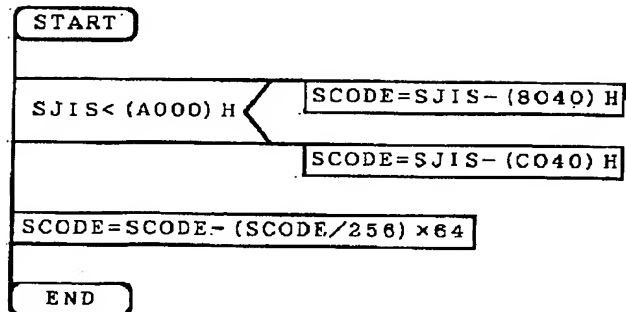
【図4】

＜文字種別型・重複排除型凝縮本文の作成方法＞



【図10】

【図10】  
＜文字コード変換アルゴリズム＞



【図21】

【図21】

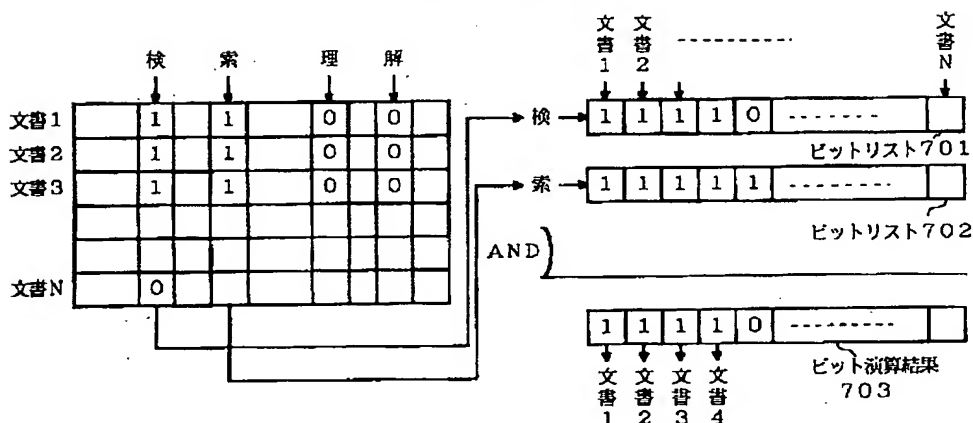
基本単語辞書

＜動詞－ある連体形＞	：ある
＜動詞－ある仮定形＞	：あれ
＜動詞－なる未然形＞	：なら
＜動詞－なる連用形＞	：なり
＜動詞－もつ未然形＞	：もた
＜動詞－もつ連体形＞	：もつ
＜助詞－が＞	：が
＜名詞－こと＞	：こと
＜名詞－ため＞	：ため
＜名詞－の＞	：の

【図7】

【図7】

（文字成分表サーチ方法）



【図6】

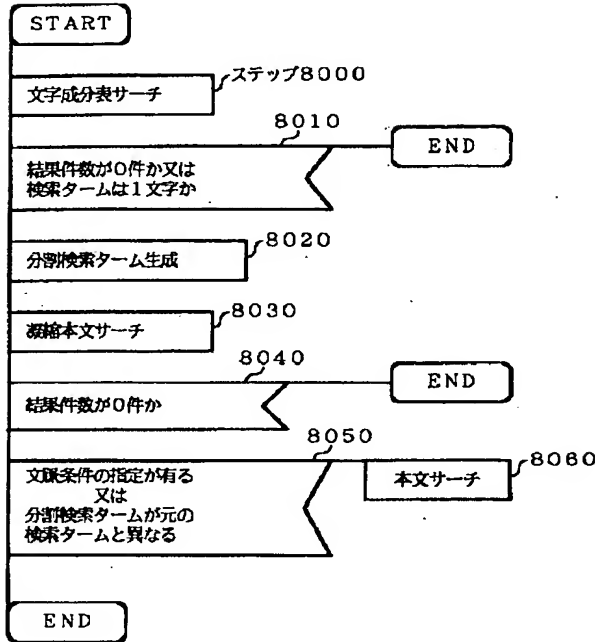
【図6】

〈文字コード依存型文字成分表の構造〉

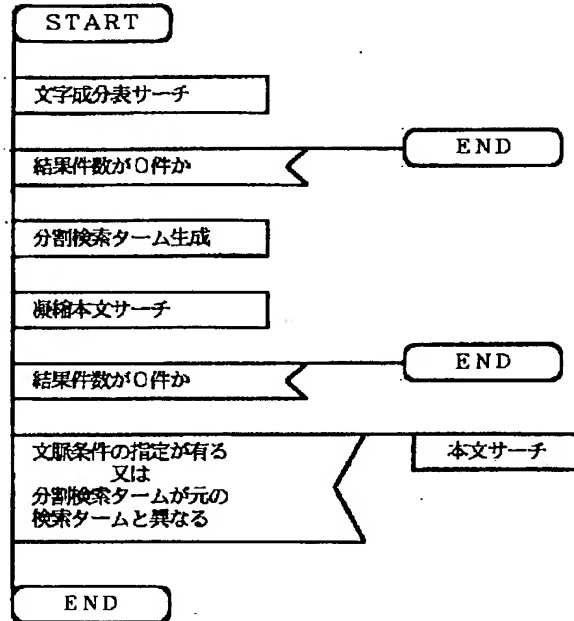
文字コード	(0000)H	(8140)H (スペース)	(82A0)H あ	(82A2)H い	(82DC)H ま	(82F0)H を	(889F)H 並	(8B5A)H 技	(8C9F)H 校	(8DF5)H 索	(E49E)H 會
文書1	0 0	....	1	1	1	....	0	....	1	....	0
文書2	0 0	....	1	0	....	....	0	....	1	....	0
文書3	0 0	....	1	0	....	....	0	....	1	....	0
文書4	0 0	....	1	0	....	1	0	0	1	....	0
.	0 0	....	1	0	....	1	0	0	1	....	0
.											
文書N											

エントリ番号 0 1 .... 33088 .... 33440 33442 .. 33500 .. 33520 .. 34975 .. 35674 .. 35999 .. 36341 .. 60062

【図8】

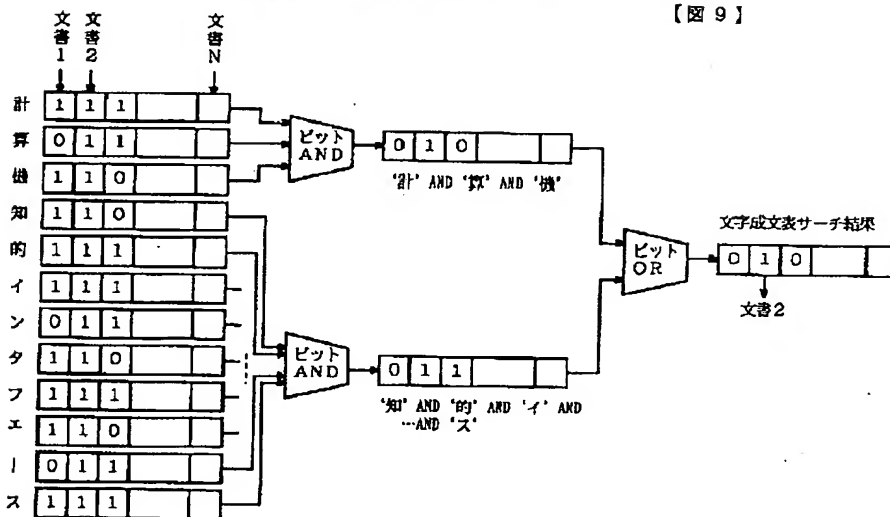
【図8】  
〈階層型プリサーチ手順〉

【図14】

【図14】  
〈ハッシング型文字成分表を用いた階層型プリサーチ手順〉

【図9】

〈複数検索タームに対する文字成分表サーチ方法〉



【図9】

【図18】

【図18】

〈ハッシュエントリ変換表〉

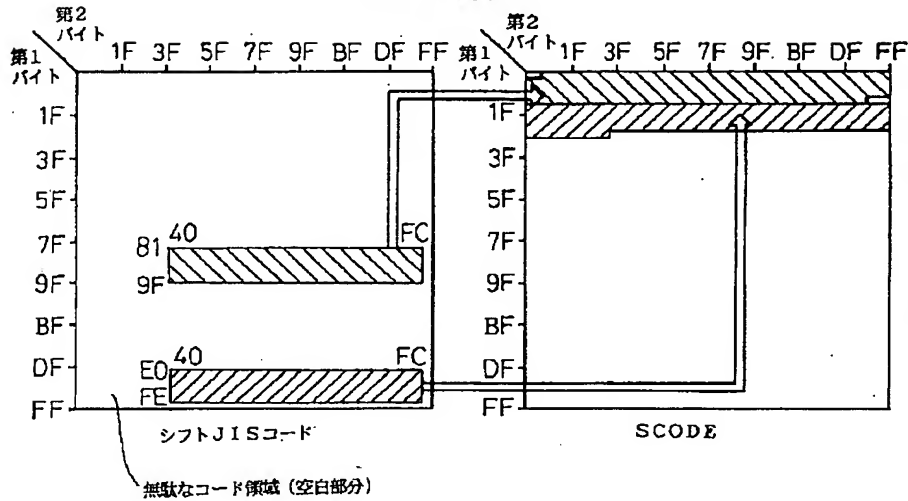
SCODE	エントリ番号
.	.
.	.
.	.
(095D)H	150
(095E)H	356
(095F)H	231
(0960)H	483
(0961)H	2
(0962)H	256
(0963)H	25
(0964)H	67
.	.
.	.
.	.

【図11】

【図33】

【図11】

&lt;文字コード変換方法&gt;



【図33】

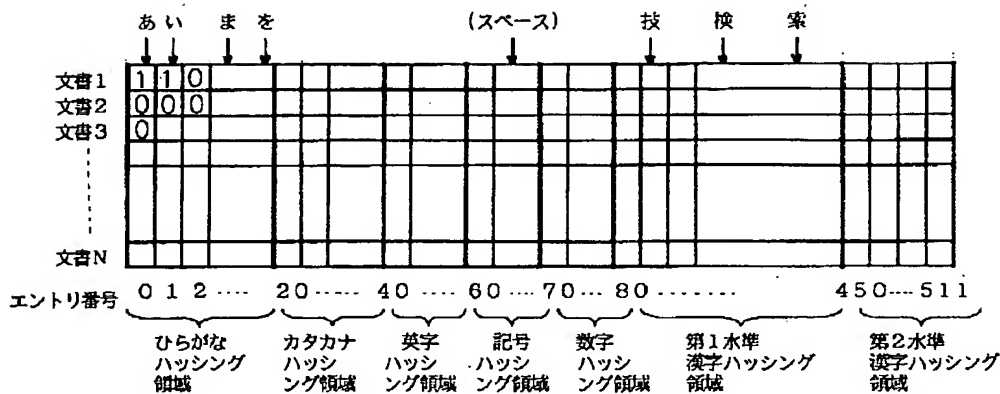
&lt;文字として格納した文字成分表&gt;

文書1	あいま検索のため技術・・・
文書2	自然語による検索技術・・・
文書3	壁を検出しながら出口探索・・・
文書4	文書理解を用いた検索システム・・・
・	・・・・・・
・	・・・
文書N	

【図15】

【図15】

&lt;文字種を考慮したハッシングによる文字成分表の構造&gt;





【図12】

【図12】

〈文字コード変換を用いた文字成分表の構造〉

変換文字コード	(0000)H	(スペース) (00C0)H	あ (01E0)H	い (01E2)H	ま (021C)H	を (0230)H	皿 (065F)H	技 (085A)H	校 (095F)H	索 (0A75)H	位 (1FDE)H
文書1	0 0	....	1	....	1	....	0	....	1	....	1 ...
文書2	0 0	....	1	....	0	....	0	....	1	....	1 ...
文書3	0 0	....	1	....	0	....	0	....	1	....	1 ...
文書4	0 0	....	1	....	0	....	0	....	1	....	1 ...
.	0 0	....	1	....	0	....	0	....	1	....	1 ...
.	0 0	....	1	....	0	....	0	....	1	....	1 ...
文書N											

エントリ番号 0 1 ..... 192 .... 480 482 .. 540 .. 560 .. 1631 .. 2138 .. 2399 .. 2677 ... 8158

【図13】

【図13】

＜ハッシュ化エントリによる文字成分表の構造＞

$h(SCODE) = \text{mod}(SCODE, 512)$  の場合

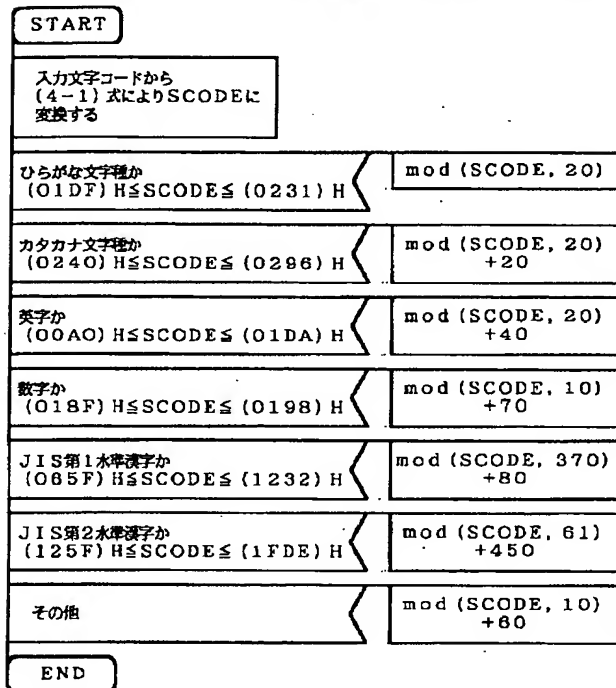
ハッシュ化 文字コード	を 技 画 梁 (スペース) 換 あ い															
	(0)	ま (28)	を (48)	を (90)	を (95)	を (117)	(192)	(351)	あ (480)	い (482)						
文書1	0	0	...	1	...	0	...	1	...	1	0	1	...	0	...	0
文書2	0	0	...	1	...	0	...	1	...	1	0	0	...	0	...	0
文書3	0	0	...	0	...	0	...	1	...	1	0	0	...	0	...	0
文書4	0	0	...	0	...	0	...	1	...	1	0	0	...	0	...	0
.	0	0	...	0	...	0	...	1	...	1	0	0	...	0	...	0
.	0	0	...	0	...	0	...	1	...	1	0	0	...	0	...	0
文書N																

エントリ番号 0 1 ... 28 ... 48 ... 90 .. 95 ... 117 .... 192 .... 351 ... 480 482 ... 511

【図16】

【図16】

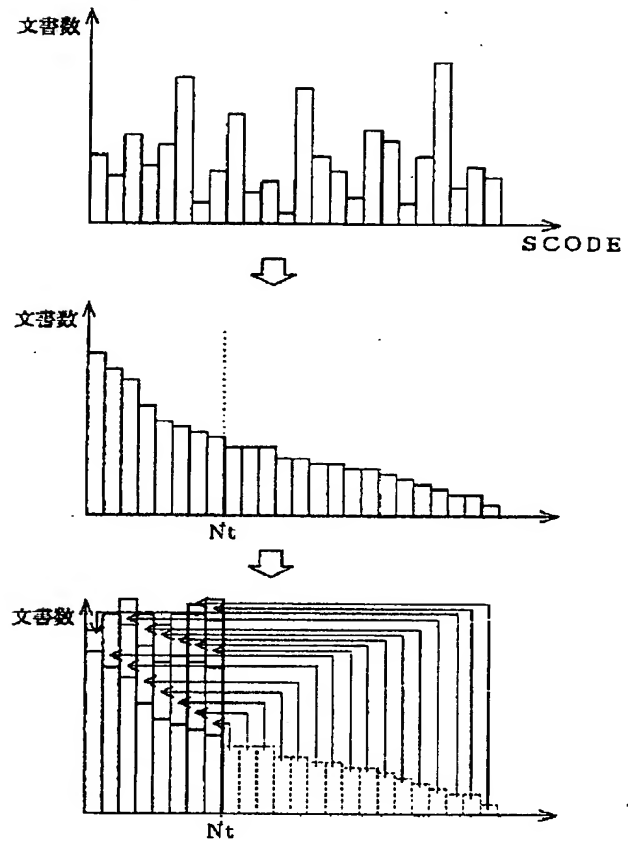
〈文字種を考慮したハッシングによる文字成分表を用いた階層型ブリスーチ手順〉



【図17】

【図17】

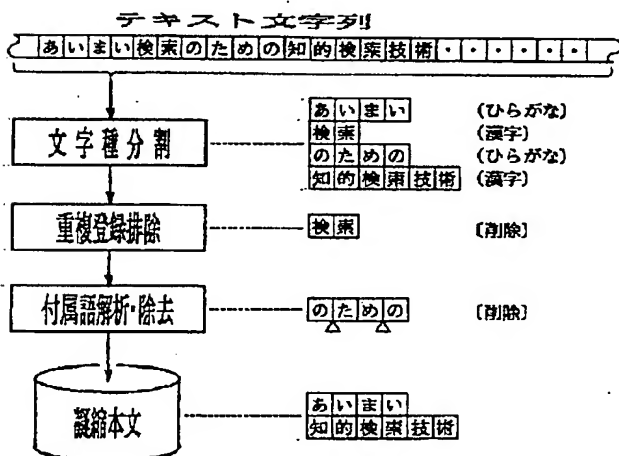
〈頻度情報を考慮したハッシング方法〉



【図19】

【図19】

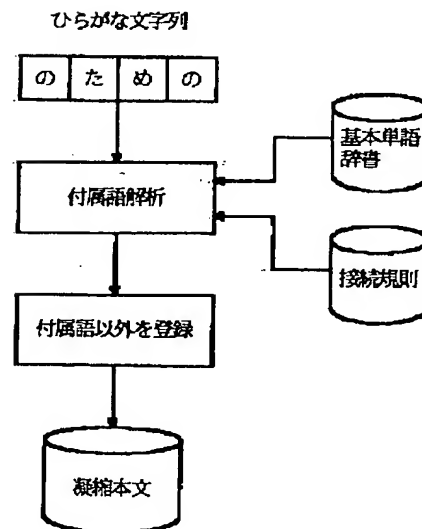
〈文字種分割・重複語排除・付属語除去型凝縮本文の作成方法〉



【図20】

【図20】

〈付属語除去方法〉



【図22】

【図22】

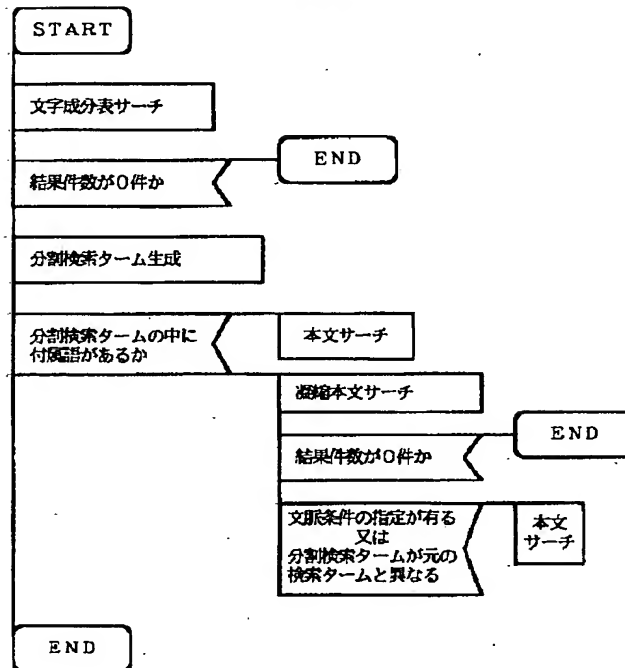
## 接続規則

- 接続1: <動詞-もつ連体形> + <名詞-こと>  
 接続2: <動詞-もつ連体形> + <名詞-ため>  
 接続3: <名詞-こと> + <助詞-が>  
 接続5: <動詞-する連体形> + <名詞-こと>  
 接続6: <助詞-の> + <名詞-ため>  
 接続7: <名詞-ため> + <助詞-の>  
 接続8: <接続詞> + <助詞-は>  
 接続9: <動詞-する未然形> + <助詞-ば>  
 接続10: <名詞-こと> + <助詞-に>  
 接続11: <助詞-で> + <助詞-ある>  
 接続12: <動詞-ある未然形> + <助詞-ば>  
 接続13: <動詞-ある連用形> + <助詞-た>

【図23】

【図23】

<文字種分割、重複排除、付属語除去型凝縮本文を用いた階層型プリサーチ手順>

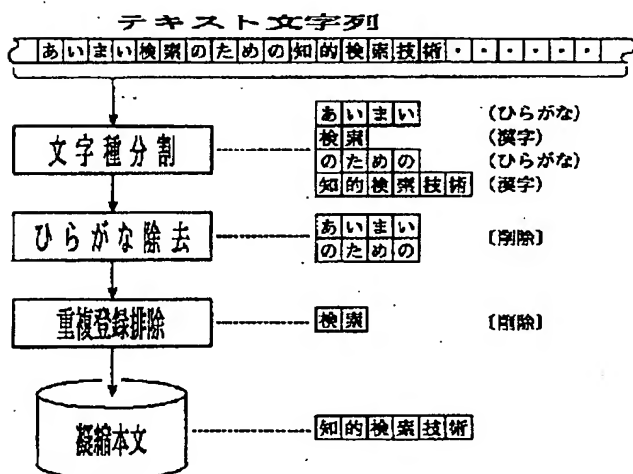


【図24】

【図25】

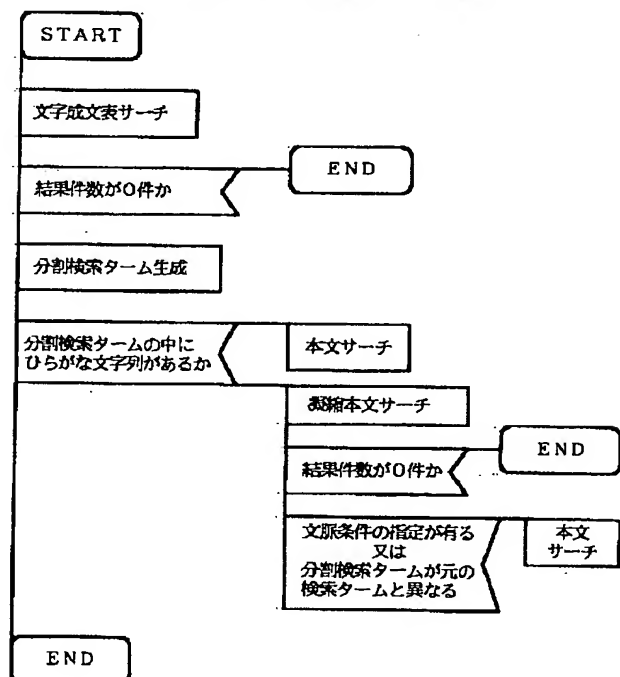
【図24】

<文字種分割・重複排除・ひらがな除去型凝縮本文の作成方法>



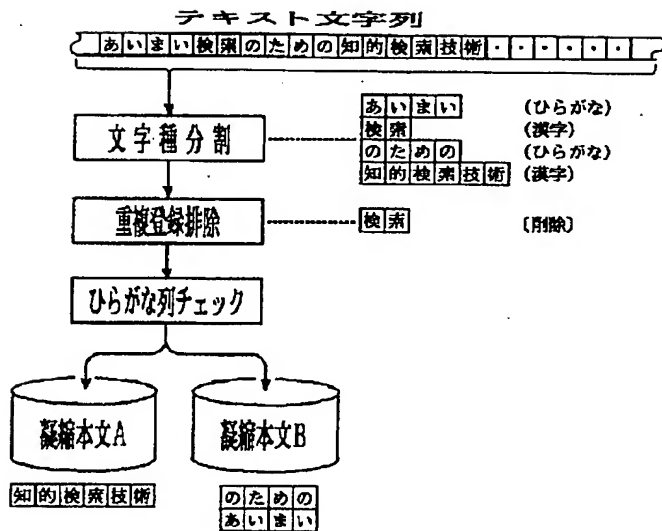
【図25】

<文字種分割、重複排除、ひらがな除去型凝縮本文を用いた階層型プリサーチ手順>



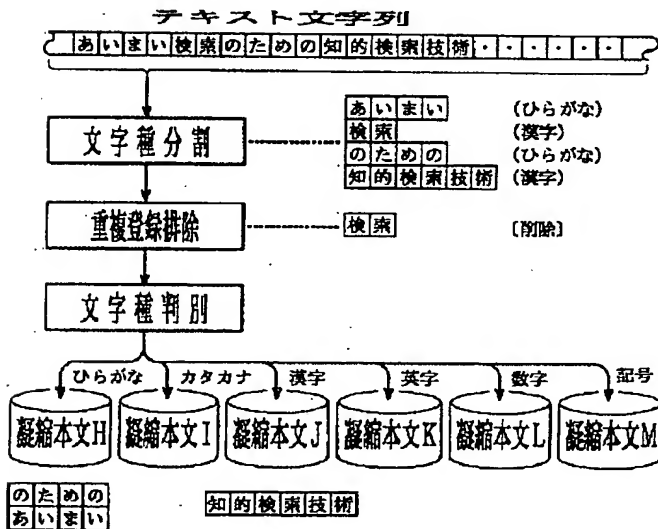
【図 26】

【図 26】



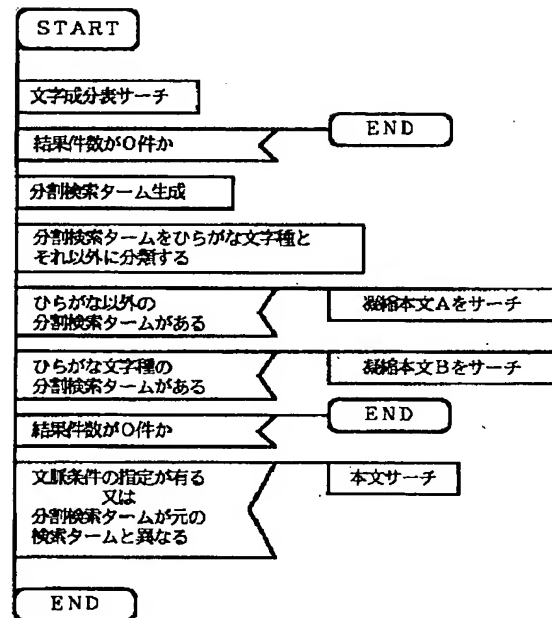
【図 28】

【図 28】



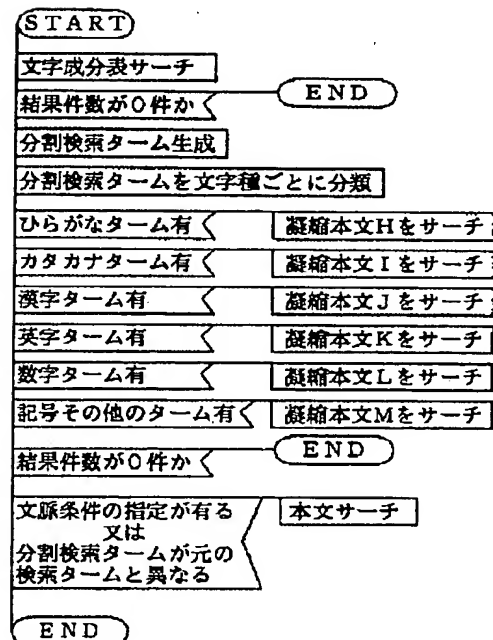
【図 27】

【図 27】

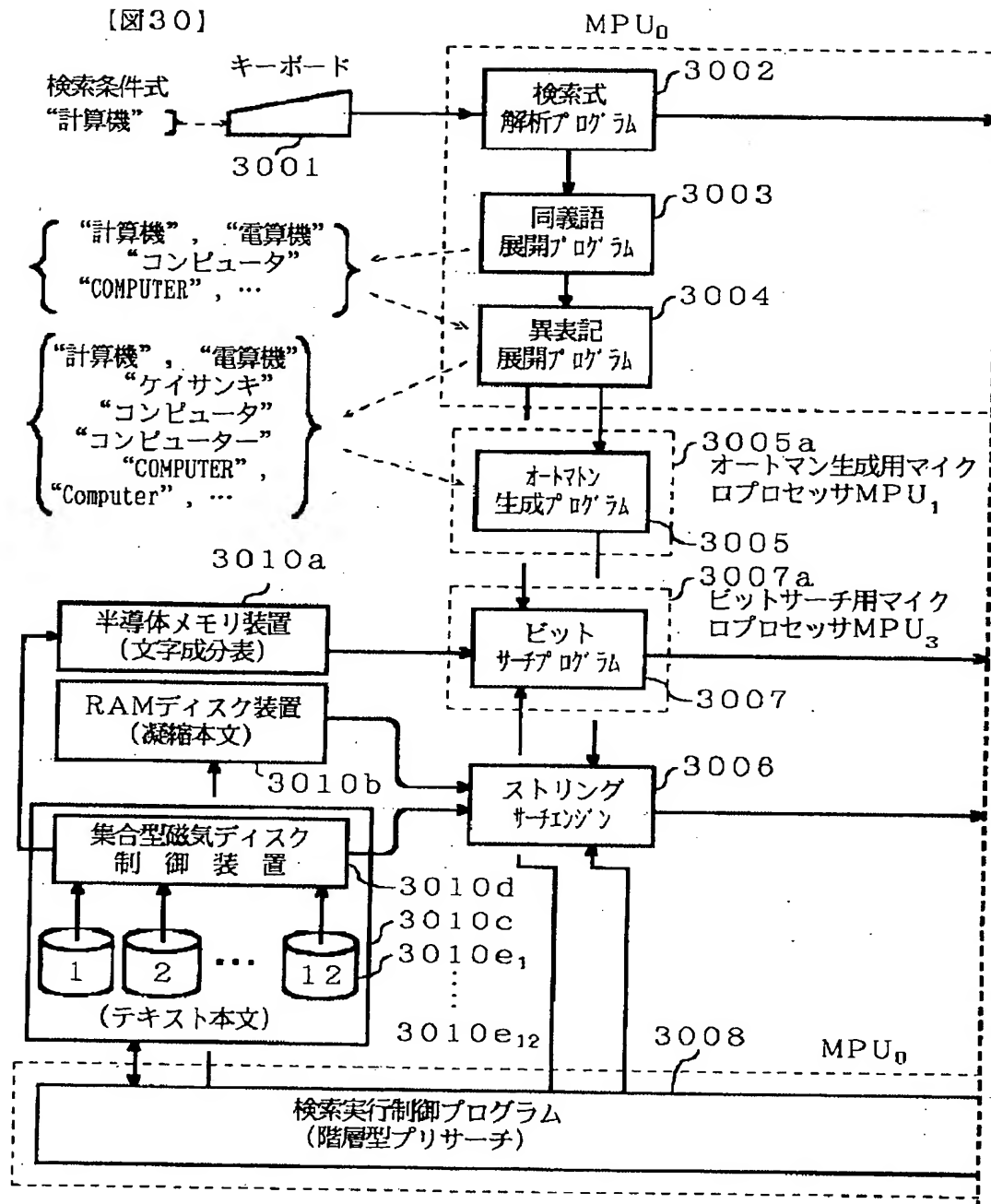


【図 29】

【図 29】

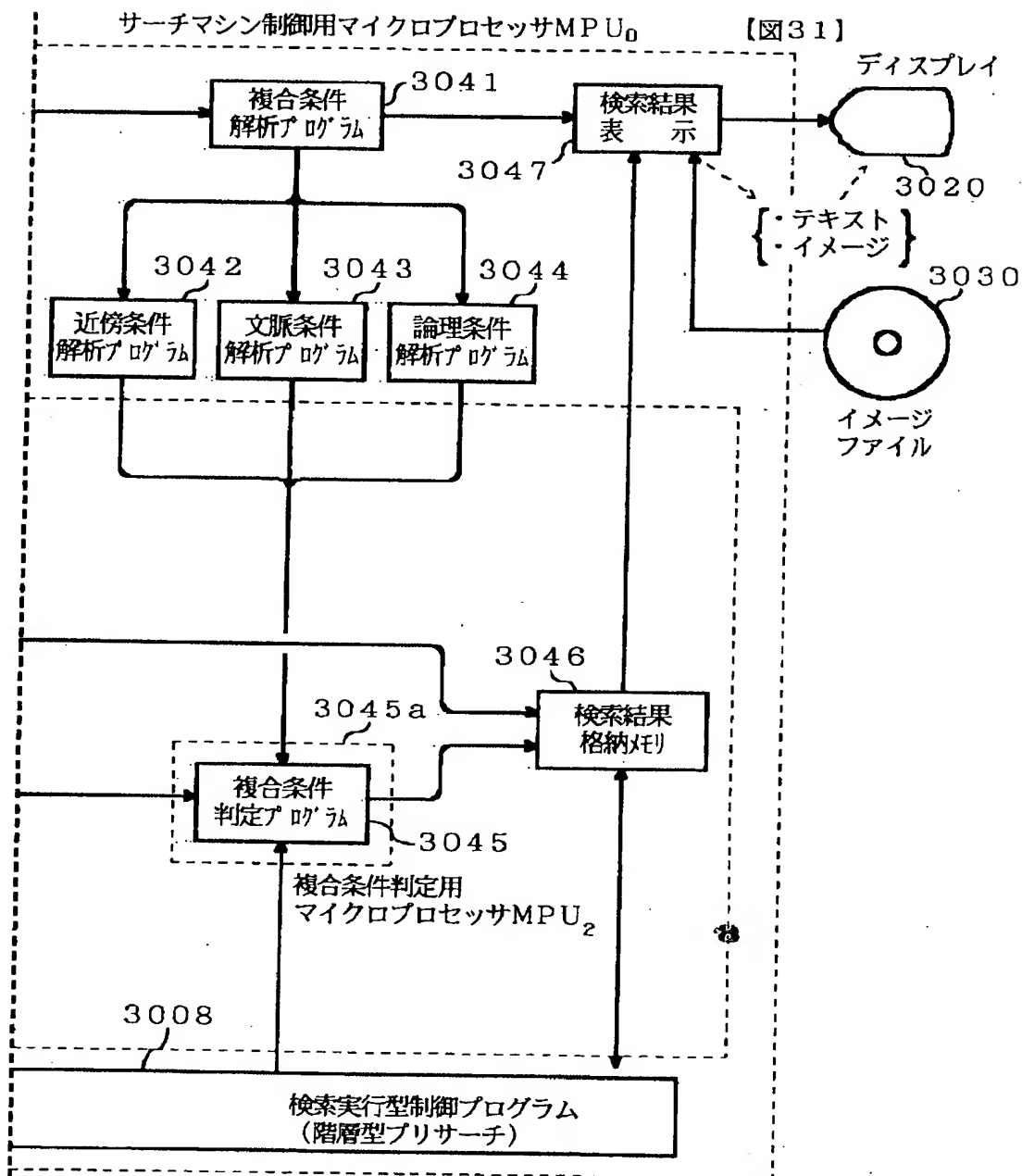


【図30】





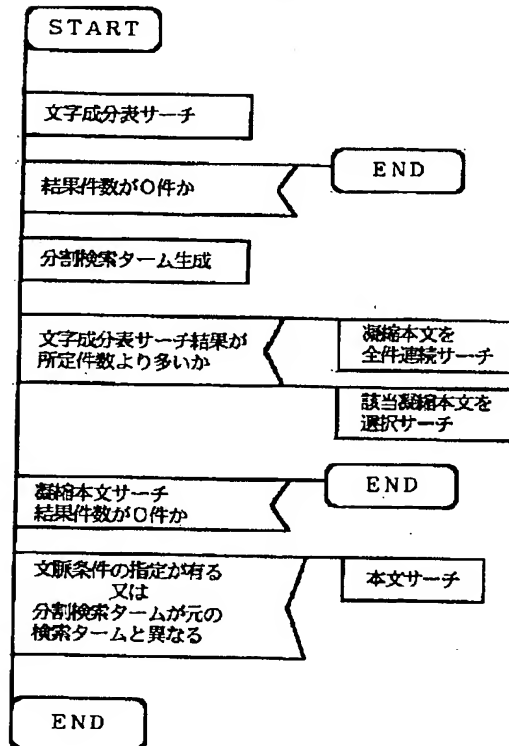
【图31】



【図 32】

【図32】

〈凝縮本文を磁気ディスクへ格納した場合の  
随層ブリサーチ手順〉



フロントページの続き

(72)発明者 加藤 寛次  
東京都国分寺市東恋ヶ窪 1 丁目280番地  
株式会社日立製作所中央研究所内  
(72)発明者 川口 久光  
東京都国分寺市東恋ヶ窪 1 丁目280番地  
株式会社日立製作所中央研究所内

(72)発明者 嶺岸 直材  
神奈川県川崎市幸区鹿島田890番地の12  
株式会社日立製作所情報システム工場内  
Fターム(参考) 5B075 ND03 NR02 NR14 NS10 PP02  
PQ02 QM01